
Sparse Kernel Machines

*Prof. Nicholas Zabaras
Center for Informatics and Computational Science*

<https://cics.nd.edu/>

*University of Notre Dame
Notre Dame, Indiana, USA*

Email: nzabaras@gmail.com

URL: <https://www.zabaras.com/>

April 3, 2019

SVM - Contents

- ❑ [Support Vector Machines](#), [Linearly Separable Training Data](#), [Maximum Margin Classifiers](#), [Constrained Optimization Problem](#), [Dual Formulation](#), [Quadratic Programming](#), [KKT Conditions](#), [Equivalent Error Function Formulation](#), [Example – SVM with a Gaussian Kernel](#)
- ❑ [Overlapping Class Distributions](#), [v-SVM](#), [Chunking](#), [Decomposition Methods](#), [Probabilistic Predictions with SVM](#)
- ❑ [SVM and Logistic Regression](#), [Multiclass SVMs](#)
- ❑ [SVM for Regression](#), [Computational Learning Theory](#)

Following closely:

Bishop CM, [Pattern Recognition and Machine Learning](#), Springer, 2006 (Chapter 7)

Murphy, K. [Machine Learning: A probabilistic Approach](#) (Chapter 14)

Machine Learning, University of Notre Dame, Notre Dame, IN, USA (Spring 2019, N. Zabarás)

RVM - Contents

- ❑ Relevance Vector Machines: Likelihood, Prior, Posterior, Predictive Distribution
- ❑ Evidence Approximation, Evidence, Two useful matrix identities
- ❑ Fast training algorithm
- ❑ Splitting the evidence, Analysis of the evidence
- ❑ Complete algorithm, Example
- ❑ Fast updates
- ❑ Extension to Multiple Outputs, Numerical Stability Concerns

Support Vector Machines

Introduction

- ❑ Significant limitation of earlier discussed algorithms is that $k(\mathbf{x}_n, \mathbf{x}_m)$ must be evaluated for all possible $\mathbf{x}_n, \mathbf{x}_m$ of training points
 - ✓ Computationally infeasible during training
 - ✓ Excessive computation when making predictions for test data.
- ❑ In this lecture, we discuss kernel-based algorithms that have *sparse solutions*
 - ✓ predictions for new inputs depend only on the kernel function evaluated at a subset of the training data points.
- ❑ In the *support vector machine* (SVM), the determination of the model parameters corresponds to a *convex optimization problem*, i.e. any local solution is also a global optimum.
- ❑ However, SVM does not provide probabilistic solutions.

Linearly Separable Training Data

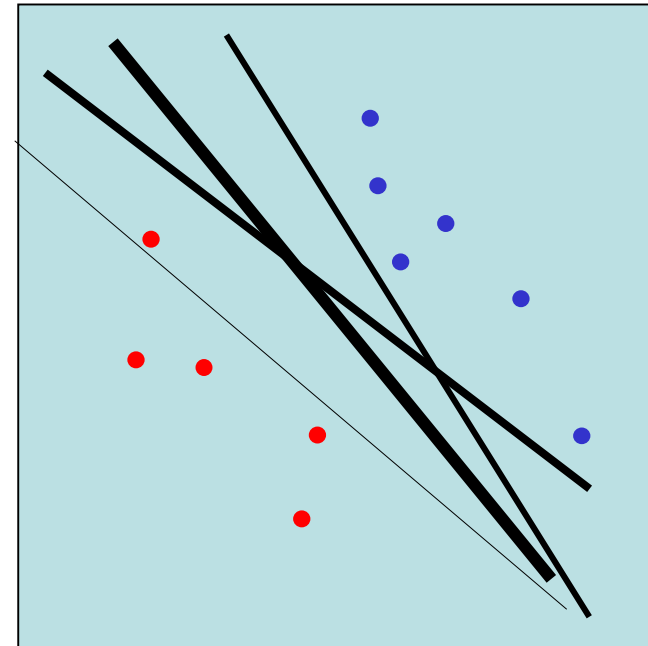
- Consider two-class classification problem using linear models

$$y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$$

- The training data set comprises N input vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$, with corresponding target values t_1, \dots, t_N where $t_n \in \{-1, 1\}$, and new data points \mathbf{x} are classified according to the sign of $y(\mathbf{x})$.
- Assume that the training data set is linearly separable in feature space: there exists at least one choice of \mathbf{w} and b such that $y(\mathbf{x}_n) > 0$ for points having $t_n = +1$ and $y(\mathbf{x}_n) < 0$ for points having $t_n = -1$, so that $t_n y(\mathbf{x}_n) > 0$ for all training data points.
- Vapnik, V. N. (1995). *The nature of statistical learning theory*. Springer.
- Burges, C. J. C. (1998). [A tutorial on support vector machines for pattern recognition](#). *Knowledge Discovery and Data Mining* 2(2), 121–167.
- Cristianini, N. and J. Shawe-Taylor (2000). *Support vector machines and other kernel-based learning methods*. Cambridge University Press.
- Müller, K. R., S. Mika, G. Rätsch, K. Tsuda, and B. Schölkopf (2001). [An introduction to kernel based learning algorithms](#). *IEEE Transactions on Neural Networks* 12(2), 181–202.
- Schölkopf, B. and A. J. Smola (2002). *Learning with Kernels*. MIT Press.
- Herbrich, R. (2002). *Learning Kernel Classifiers*. MIT Press.

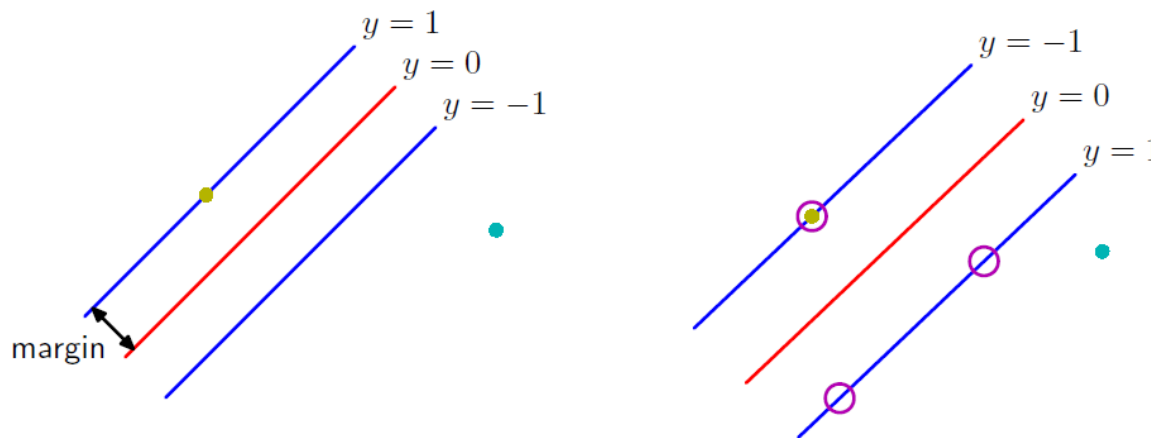
Linearly Separable Training Data

- Suppose we use a big set of features to ensure that the two classes are linearly separable. What is the best separating line to use?
- ✓ The Bayesian answer is to use them all (including ones that do not quite separate the data.)
- ✓ Weight each line by its posterior probability (i.e. by a combination of how well it fits the data and how well it fits the prior).
- ✓ Is there an efficient way to approximate the correct Bayesian answer?



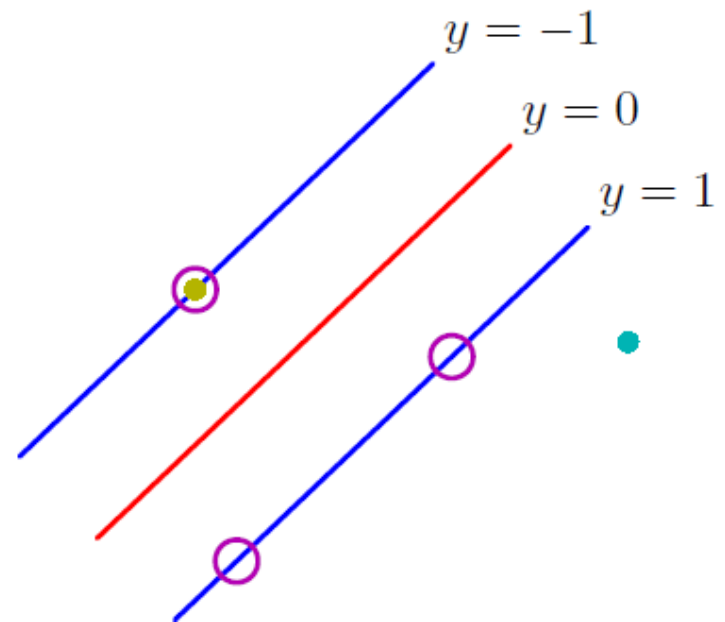
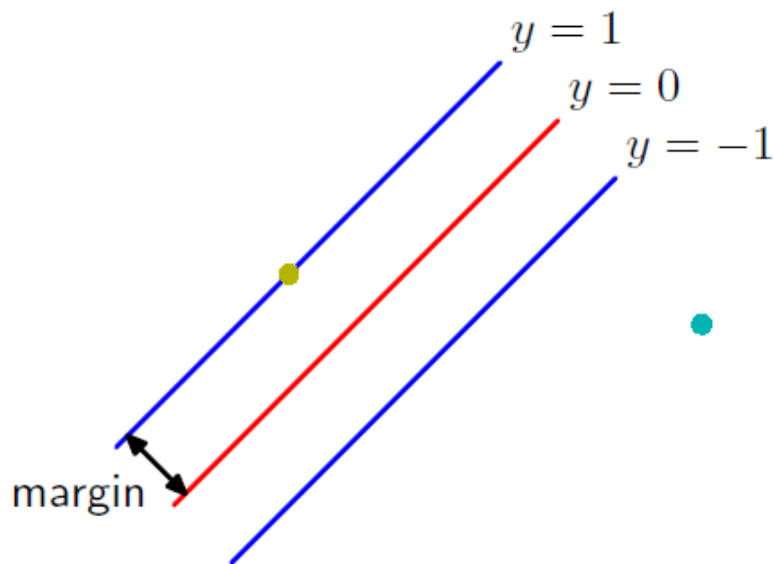
Maximum Margin Classifiers

- There may exist many solutions that separate the classes exactly.
- E.g. the perceptron algorithm finds a solution in a finite number of steps that depends on the initial values of w and b and the order in which the data points are presented.
- When multiple solutions find the one that gives the smallest generalization error.
- SVM approaches this problem through the *margin* - smallest distance between the decision boundary and any of the samples.

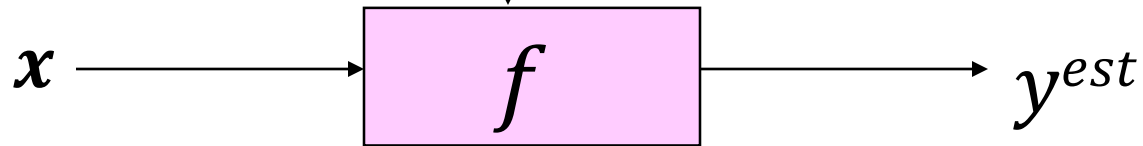


Maximum Margin Classifiers

- ❑ The margin is defined as the perpendicular distance between the decision boundary and the closest of the data points.
- ❑ Maximizing the margin leads to a particular choice of decision boundary.
- ❑ The location of this boundary is determined by a subset of the data points, known as **support vectors** (circles).

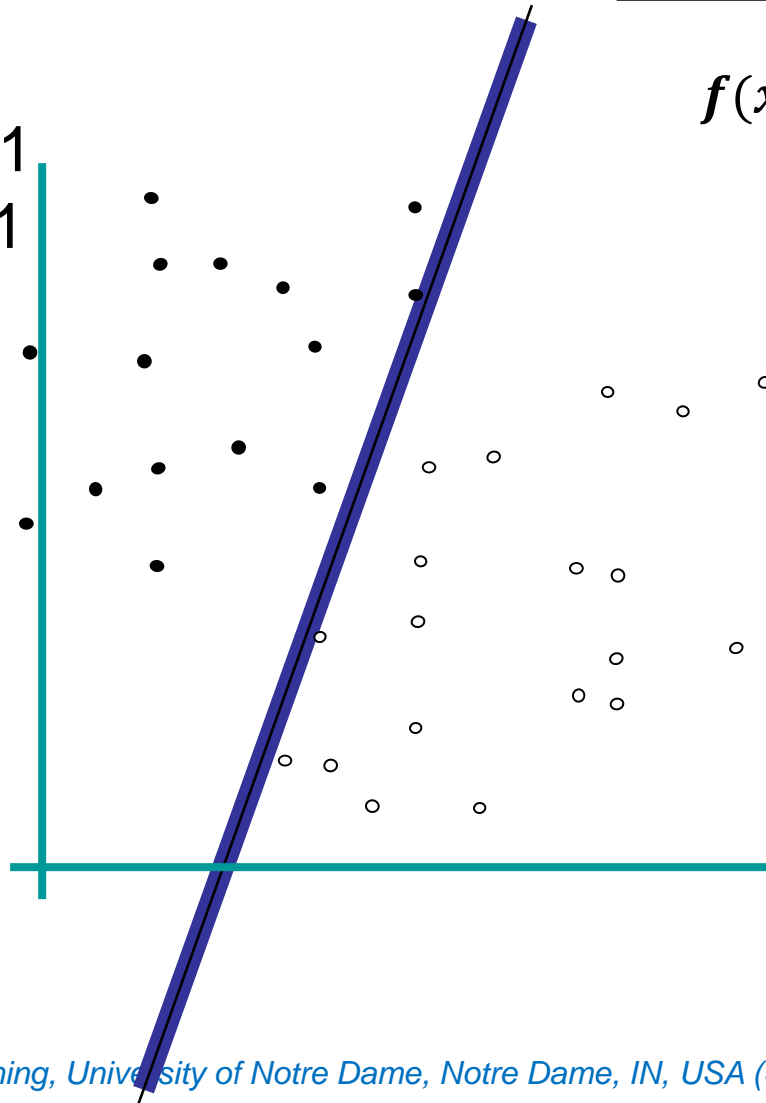


Classifier Margin



$$f(x, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot x + b)$$

- denotes +1
- denotes -1

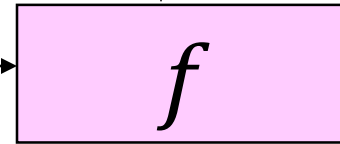


Linear SVM

Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

Support Vectors

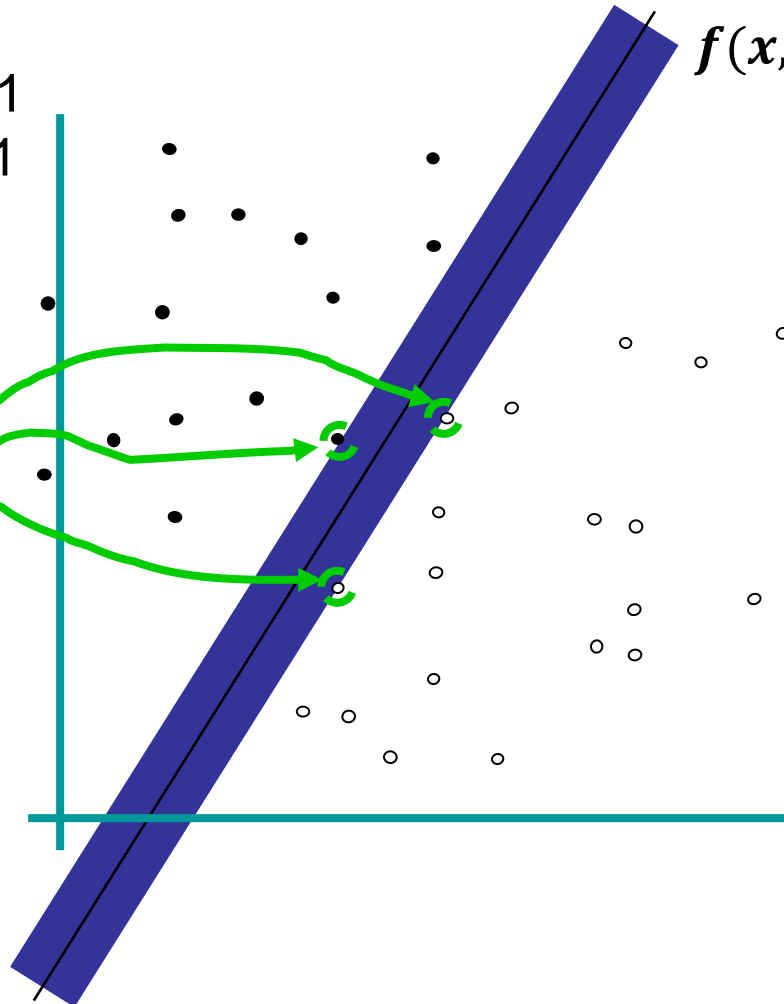
x



y^{est}

- denotes +1
- denotes -1

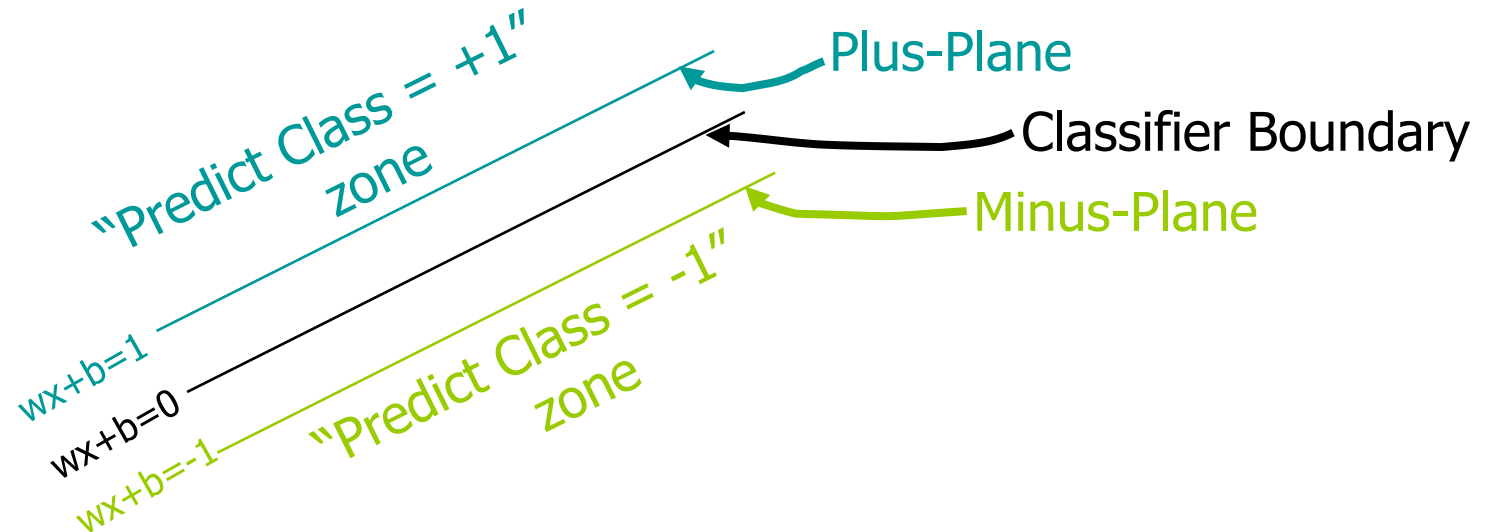
Support Vectors
are those
datapoints that
the margin
pushes up
against



$$f(x, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot x + b)$$

The maximum
margin linear
classifier is the
linear classifier
with the maximum
margin.

Specifying a line and margin



□ Plus-plane = $\{x : w \cdot x + b = +1\}$

□ Minus-plane = $\{x : w \cdot x + b = -1\}$

Classify as..

+1	if	$w \cdot x + b \geq 1$
-1	if	$w \cdot x + b \leq -1$

Maximum Margin Classifiers

- ❑ The maximum margin can be motivated with *statistical learning theory*.
- ❑ An insight has been given by Tong and Koller (2000) who consider a hybrid (generative and discriminative) classification approach.
- ❑ They first model the distribution over x for each class using a Parzen density estimator with Gaussian kernels having a common σ^2 . Together with the class priors, this defines an optimal misclassification-rate decision boundary.
- ❑ Instead of using this optimal boundary, they determine the best hyperplane by minimizing the probability of error relative to the learned density model. As $\sigma^2 \rightarrow 0$, the optimal hyperplane becomes one having maximum margin.
- ❑ As $\sigma^2 \rightarrow 0$, the hyperplane becomes independent of data points that are not support vectors.
- Tong, S. and D. Koller (2000). [Restricted Bayes optimal classifiers](#). In *Proceedings 17th National Conference on Artificial Intelligence*, pp. 658– 664. AAAI.

Maximum Margin Solution

- Since we are only interested in correctly classified points, $t_n y(\mathbf{x}_n) > 0$, the vertical distance of \mathbf{x}_n to the decision boundary $y(\mathbf{x}) = 0$ is:

$$\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

- The max margin solution is obtained as:

$$\operatorname{argmax}_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \right\}$$

- Scale \mathbf{w} and b such that for the point that is closest to the surface:

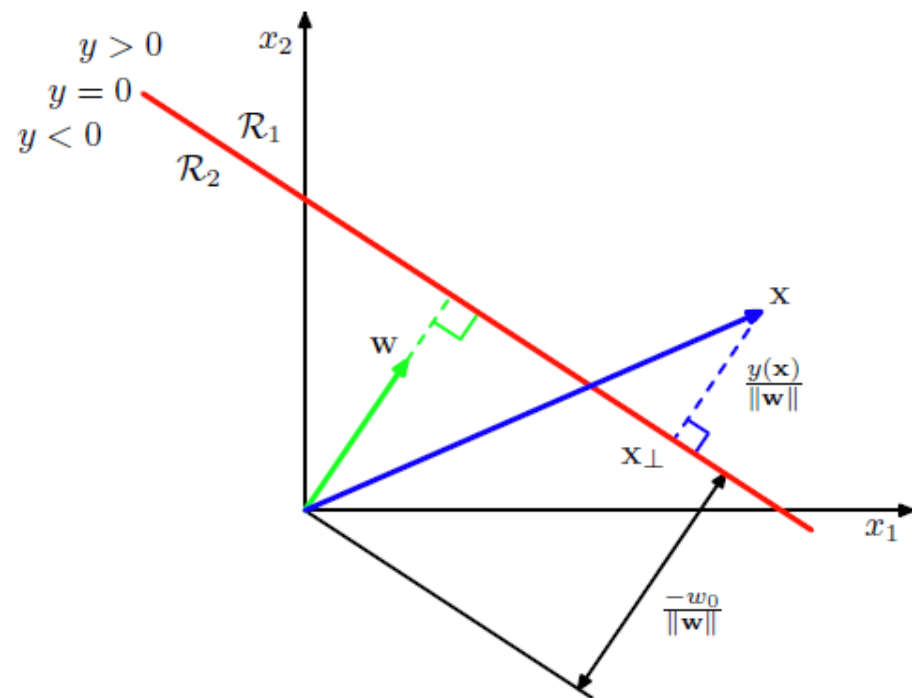
$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1$$

- Thus for all points $n = 1, \dots, N$

$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$$

- The problem now is: $\operatorname{argmin}_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$

$$\text{under } t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$$



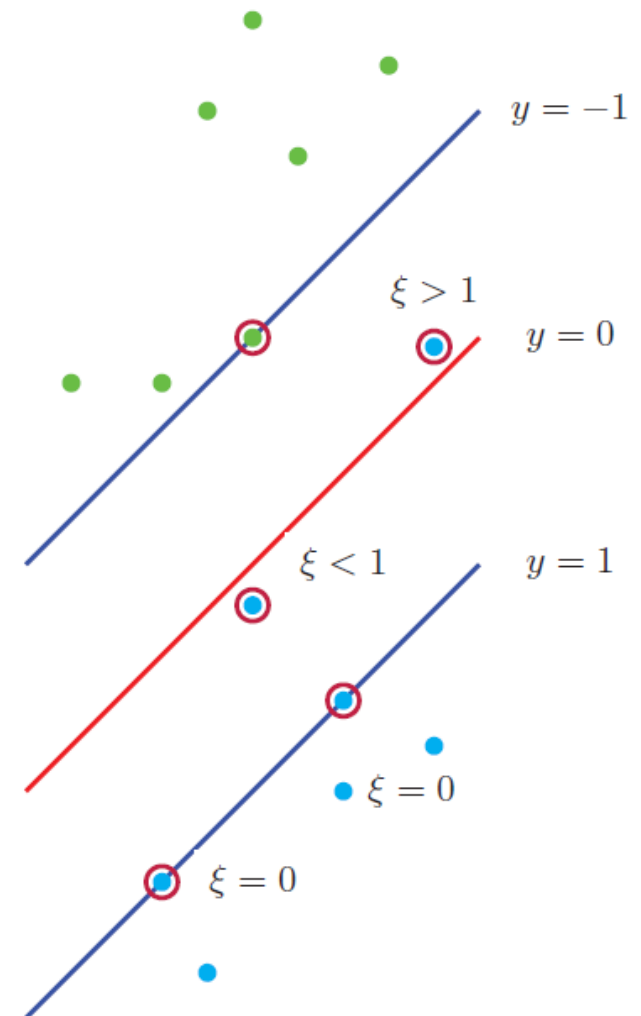
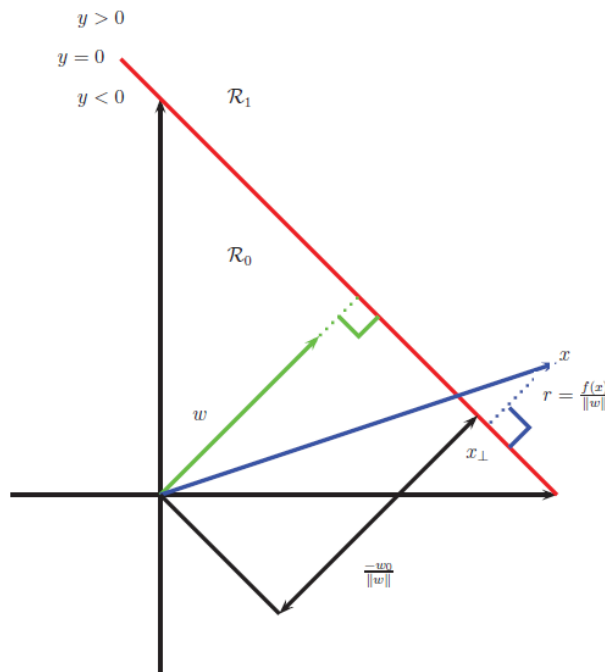
Maximum Margin Solution

□ Our problem is stated as follows:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \text{ subject to}$$

$$t_i [b + \mathbf{w}^T \phi(\mathbf{x}_i)] - 1 \geq 0, i = 1, \dots, N$$

(each point is on the correct side of the boundary)



Constrained Optimization

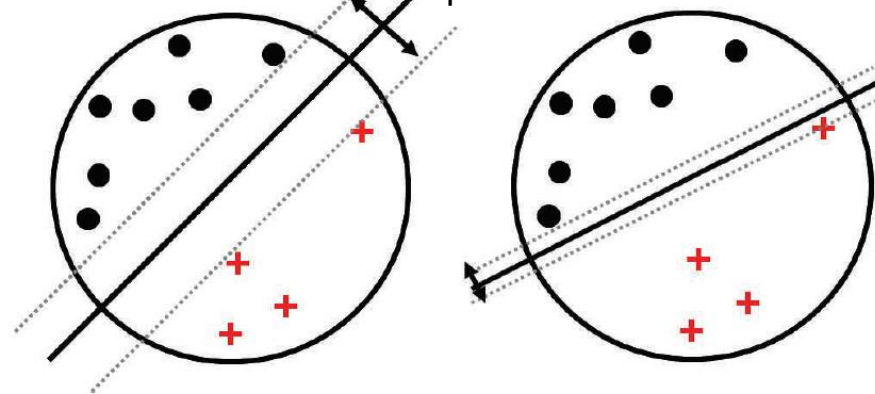
- Our problem is stated as follows:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \text{ subject to}$$

$$t_i [b + \mathbf{w}^T \phi(\mathbf{x}_i)] - 1 \geq 0, i = 1, \dots, N$$

(each point is on the correct side of the boundary)

Margin here defined as the perpendicular distance to the closest point



- We represent the constraints with Lagrange multipliers as losses:

$$\max_{\alpha \geq 0} \alpha_i (1 - t_i [b + \mathbf{w}^T \phi(\mathbf{x}_i)]) = \begin{cases} 0, & \text{if } t_i [b + \mathbf{w}^T \phi(\mathbf{x}_i)] - 1 \geq 0 \\ \infty, & \text{otherwise} \end{cases}$$

- The minimization problem can now be stated as:

$$\min_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \max_{\alpha \geq 0} \alpha_i (1 - t_i [b + \mathbf{w}^T \phi(\mathbf{x}_i)]) \right\}$$

SVM - Separable Case

- The minimization problem can now be stated as:

$$\min_{\mathbf{w}} \max_{\{\alpha_i \geq 0\}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i (1 - t_i [b + \mathbf{w}^T \phi(\mathbf{x}_i)]) \right\} =$$
$$\max_{\{\alpha_i \geq 0\}} \min_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i (1 - t_i [b + \mathbf{w}^T \phi(\mathbf{x}_i)]) \right\}$$

- As a result we should be able to minimize wrt \mathbf{w} for any given set of the α 's the following:

$$J(\mathbf{w}; \alpha) = \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i (1 - t_i [b + \mathbf{w}^T \phi(\mathbf{x}_i)]) \right\}$$

- Setting the derivatives wrt to \mathbf{w} and b gives:

$$\mathbf{w} - \sum_{i=1}^N \alpha_i t_i \phi(\mathbf{x}_i) = 0, \quad -\sum_{i=1}^N \alpha_i t_i = 0$$

Constrained Optimization Problem

- Introducing Lagrange multipliers $a_n \geq 0$, our objective function is:

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1)$$

- Note we minimize wrt \mathbf{w} and b and maximize wrt \mathbf{a} . Setting the derivatives with respect to \mathbf{w} and b equal to zero gives:

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n), \quad 0 = \sum_{n=1}^N a_n t_n$$

- Eliminating \mathbf{w} and b from $L(\mathbf{w}, b, \mathbf{a})$ results in the dual representation:

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m),$$

$$k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$$

$$a_n \geq 0, n = 1, \dots, N, \quad 0 = \sum_{n=1}^N a_n t_n$$

Dual Formulation

$$\begin{aligned} & \max_{\mathbf{a}} \tilde{L}(\mathbf{a}), \\ \tilde{L}(\mathbf{a}) &= \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m), \\ k(\mathbf{x}_n, \mathbf{x}_m) &= \boldsymbol{\phi}(\mathbf{x}_n)^T \boldsymbol{\phi}(\mathbf{x}_m) \\ a_n &\geq 0, n = 1, \dots, N, 0 = \sum_{n=1}^N a_n t_n \end{aligned}$$

- ❑ The solution to a quadratic programming problem in M variables is $\mathcal{O}(M^3)$.
- ❑ The dual formulation involves maximizing over N variables. For $M \ll N$, the dual problem appears disadvantageous.
- ❑ However, the model is reformulated using kernels, and so the maximum margin classifier can be applied efficiently to feature spaces whose dimensionality exceeds N including infinite feature spaces.
- ❑ $k(\mathbf{x}, \mathbf{x}')$ being positive definite ensures that $\tilde{L}(\mathbf{a})$ is bounded above, thus leading to a well posed optimization.

Quadratic Programming – KKT Conditions

- In order to classify new data points using the trained model, we evaluate the sign of $y(\mathbf{x})$ defined by $y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$
- This can be expressed in terms of $\{a_n\}$ and the kernel function by substituting $\mathbf{w} = \sum_{n=1}^N a_n t_n \boldsymbol{\phi}(\mathbf{x}_n)$ to give

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$$

- Our constraint optimization problem satisfies the following Karush-Kuhn-Tucker (KKT) conditions:

$$a_n \geq 0$$

$$t_n y(\mathbf{x}_n) - 1 \geq 0$$

$$a_n (t_n y(\mathbf{x}_n) - 1) = 0$$

KKT Conditions

$$\begin{aligned}a_n &\geq 0 \\ t_n y(\mathbf{x}_n) - 1 &\geq 0 \\ a_n(t_n y(\mathbf{x}_n) - 1) &= 0\end{aligned}$$

- Thus for every \mathbf{x}_n , either $a_n = 0$ or $t_n y(\mathbf{x}_n) = 1$.
- Any data point for which $a_n = 0$ will not appear in $y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$ and hence plays no role in making predictions for new data points.
- The remaining data points are called *support vectors*, and because they satisfy $t_n y(\mathbf{x}_n) = 1$, they correspond to points that lie on the maximum margin hyperplanes in feature space.
- Once the SVM model is trained, a significant proportion of the data points can be discarded and only the support vectors retained.

Computing b

- Having found the value of a , we can determine b by noting that any support \mathbf{x}_n satisfies $t_n y(\mathbf{x}_n) = 1$.
- Using $y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$ this gives:

$$t_n \left(\sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) + b \right) = 1$$

- Here \mathcal{S} denotes the indices of the support vectors. A numerically stable solution is obtained by multiplying by t_n , use $t_n^2 = 1$, and **average over all support vectors** before solving for b to give

$$b = \frac{1}{N_{\mathcal{S}}} \sum_{n \in \mathcal{S}} \left(t_n - \sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) \right)$$

Error Function

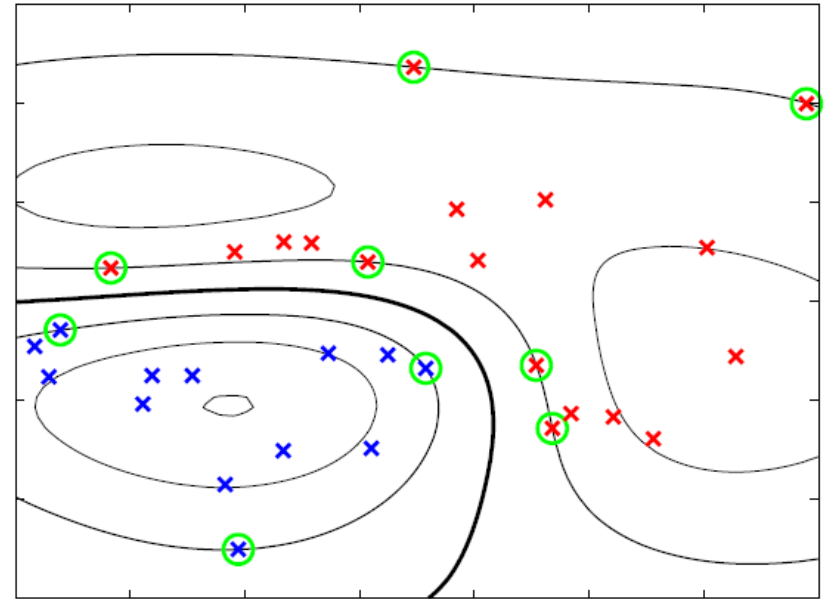
- For comparison with alternative models, we can express the maximum margin classifier in terms of minimization of an error function, with a quadratic regularizer,

$$\sum_{n=1}^N E_{\infty}(t_n y(\mathbf{x}_n) - 1) + \lambda \|\mathbf{w}\|^2$$

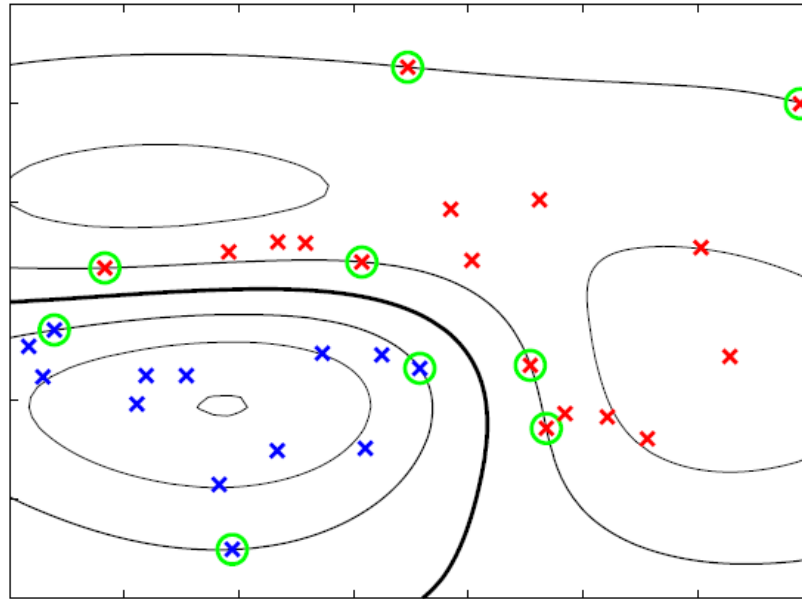
- $E_{\infty}(z)$ is a function that is zero if $z \geq 0$ and ∞ otherwise and ensures that the constraints $t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$ are satisfied.
- Note that as long as $\lambda > 0$, its precise value plays no role here.

Example: SVM with Gaussian Kernel

- Example of synthetic data from two classes in 2D showing contours of constant $y(\mathbf{x})$ obtained from a SVM having a Gaussian kernel function.
- Also shown the decision boundary, the margin boundaries, and the support vectors.
- Although the data set is not linearly separable in the 2D data space \mathbf{x} , it is linearly separable in the nonlinear feature space defined implicitly by the kernel function.
- Thus the training data points are perfectly separated in the original data space.



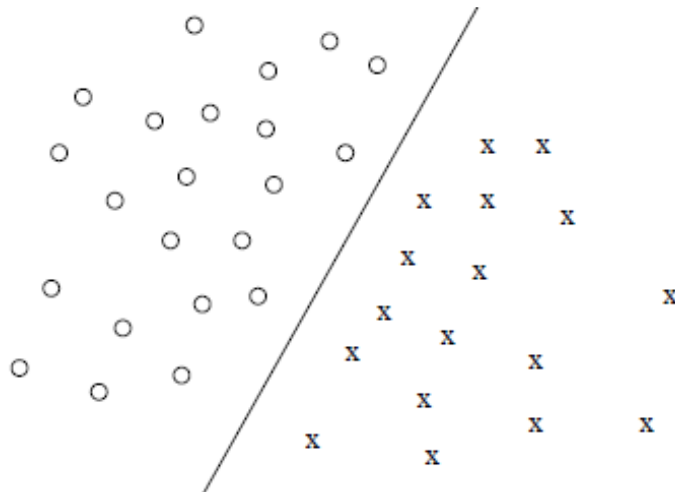
Example: SVM with Gaussian Kernel



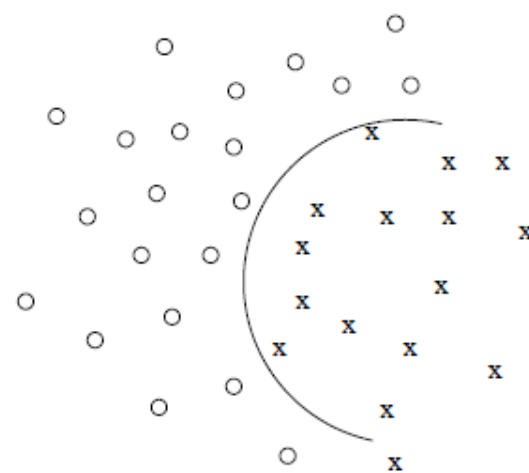
- ❑ The maximum margin hyperplane is defined by the location of the support vectors.
- ❑ Other data points can be moved around freely (so long as they remain outside the margin region) without changing the decision boundary, and so the solution remains independent of such data points.

Non-linear Classifier

- We easily obtain a non-linear classifier by using the kernel trick. Map $\mathbf{x} = [x_1 \ x_2]$ into $\boldsymbol{\phi}(\mathbf{x}) = [x_1^2 \ x_2^2 \ \sqrt{2}x_1x_2 \ \sqrt{2}x_1 \ \sqrt{2}x_2 \ 1]$



Linear separator in the
feature $\boldsymbol{\phi}$ -space

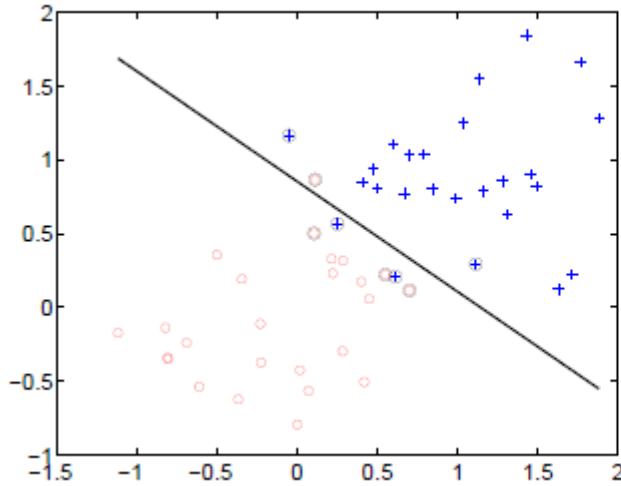


Non-linear separator
in the original \mathbf{x} -space

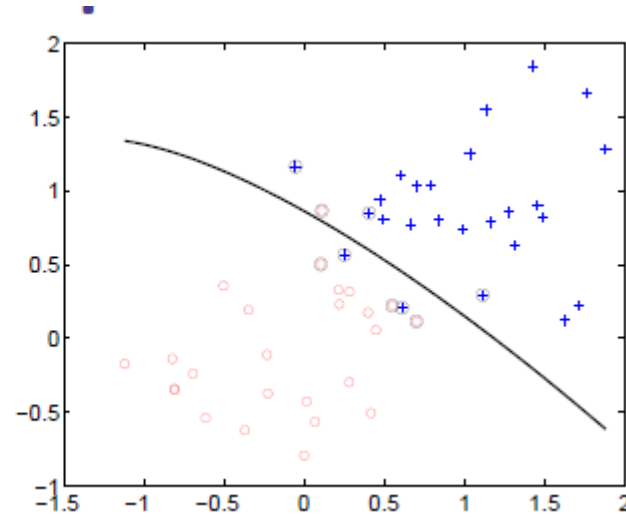
- The algorithm only needs dot products of the form:
$$\boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}') = x_1^2 x_1'^2 + x_2^2 x_2'^2 + 2x_1 x_2 x_1' x_2' + 2x_1 x_1' + 2x_2 x_2' + 1 = (1 + \mathbf{x}^T \mathbf{x}')^2$$
- So the inner products can be evaluated without ever explicitly constructing the feature vectors $\boldsymbol{\phi}(\mathbf{x})$!

SVM Classification: Examples

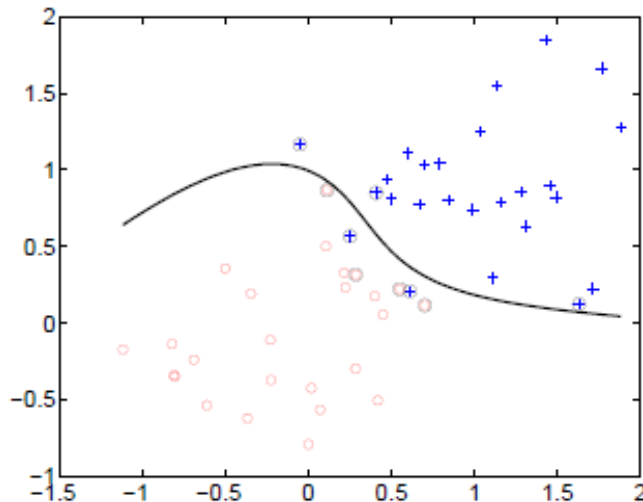
linear



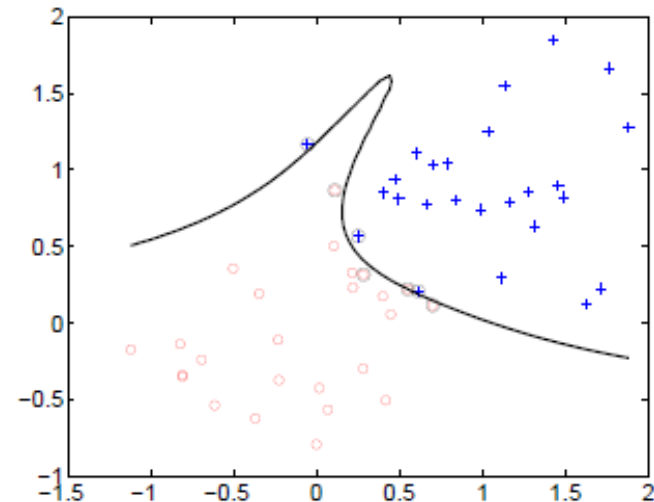
2nd order polynomial



4th order polynomial



8th order polynomial



Overlapping Class Distributions

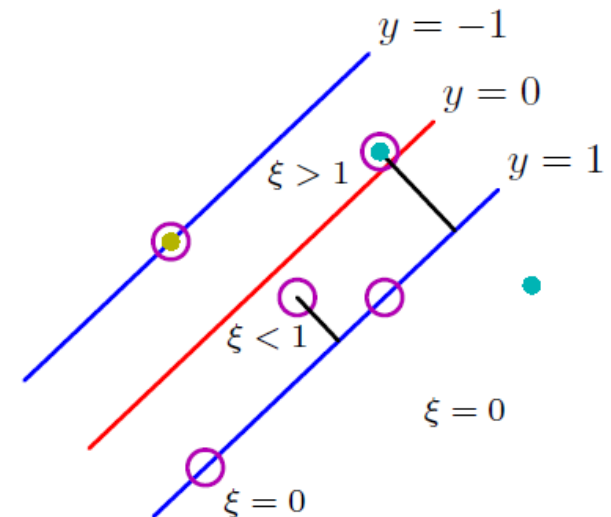
- Up to now we assumed that the training data points are linearly separable in the feature space $\phi(\mathbf{x})$. The resulting SVM gives exact separation of the training data in the space \mathbf{x} , although the corresponding decision boundary will be nonlinear.
- In practice, the class-conditional distributions may overlap, and exact separation of the training data can lead to poor generalization.
- Need to allow some of the training points to be misclassified. From

$$\sum_{n=1}^N E_{\infty}(t_n y(\mathbf{x}_n) - 1) + \lambda \|\mathbf{w}\|^2$$

we see that in the case of separable classes, we implicitly used an error function that gave infinite error if a data point was misclassified and zero error if it was classified correctly, and then optimized the model parameters to maximize the margin.

Overlapping Class Distributions

- We now modify this approach so that data points are allowed to be on the 'wrong side' of the margin boundary, but with a penalty that increases with the distance from that boundary.
- It is convenient to make this penalty a linear function of this distance.
- We introduce *slack variables*, $\xi_n \geq 0, n = 1, \dots, N$, with one slack variable for each training data point.
- $\xi_n = 0$ for data points that are on or inside the correct margin boundary
- $\xi_n = |t_n - y(x_n)|$ for other points. Thus a data point that is on the decision boundary $y(x_n) = 0$ will have $\xi_n = 1$, and points with $\xi_n > 1$ will be misclassified.



- Bennett, K. P. (1992). [Robust linear programming discrimination of two linearly separable sets](#). *Optimization Methods and Software* **1**, 23–34.
- Cortes, C. and V. N. Vapnik (1995). [Support vector networks](#). *Machine Learning* **20**, 273–297.

Overlapping Class Distributions

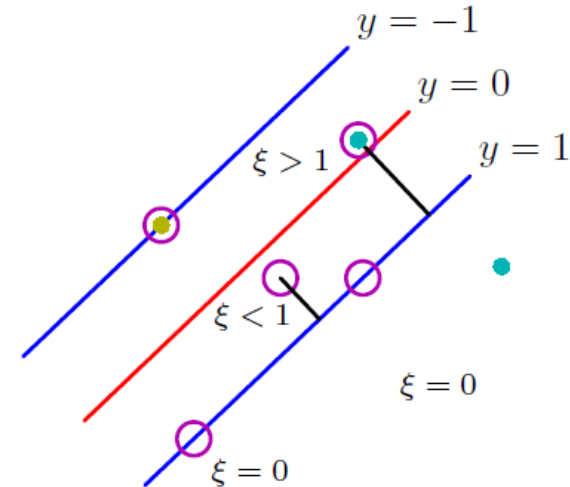
- The exact classification constraints

$$t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$$

are then replaced with

$$t_n y(\mathbf{x}_n) \geq 1 - \xi_n, n = 1, \dots, N$$

in which the slack variables are constrained to satisfy $\xi_n \geq 0$.



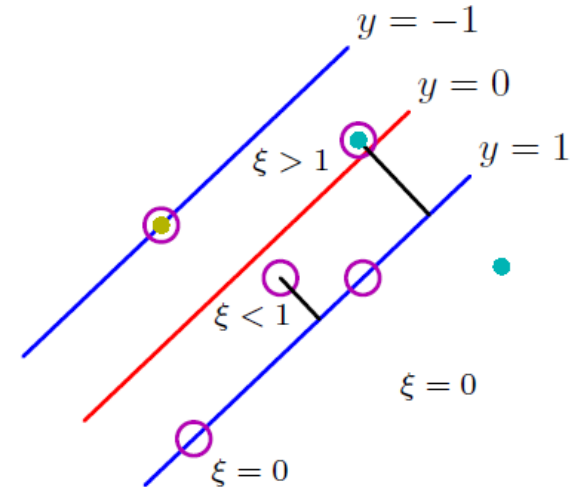
- Data points for which $\xi_n = 0$ are correctly classified and are either on the margin or on the correct side of the margin.
- Points for which $0 < \xi_n \leq 1$ lie inside the margin, but on the correct side of the decision boundary, and those data points for which $\xi_n > 1$ lie on the wrong side of the decision boundary and are misclassified.
- This is the *soft margin* constraint.

Overlapping Class Distributions

- ❑ This allows some of the training data points to be misclassified.
- ❑ While slack variables allow for overlapping class distributions, this framework is still sensitive to outliers because the penalty for misclassification increases linearly with ξ .
- ❑ Our goal is now to maximize the margin while softly penalizing points that lie on the wrong side of the margin boundary. We minimize

$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

where $C > 0$ controls the trade-off between the slack variable penalty and the margin.



Overlapping Class Distributions

$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

- Because any point that is misclassified has $\xi_n > 1$, it follows that $\sum_{n=1}^N \xi_n$ is an upper bound on the number of misclassified points.
- C is like the inverse of a regularization coefficient because it controls the trade-off between minimizing training errors and controlling model complexity.
- As $C \rightarrow \infty$, we will recover the earlier support vector machine for separable data.

Overlapping Class Distributions

- We now wish to minimize

$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

subject to the constraints

$$t_n y(\mathbf{x}_n) \geq 1 - \xi_n, n = 1, \dots, N$$

together with $\xi_n \geq 0$.

- The corresponding Lagrangian is given by

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \mathbf{a}, \boldsymbol{\mu}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n$$

where $a_n \geq 0, \mu_n \geq 0$ are Lagrange multipliers for the two constraints.

KKT Conditions

$$L(\mathbf{w}, b, \xi, \mathbf{a}, \boldsymbol{\mu}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n$$

□ The KKT conditions for $n = 1, \dots, N$ are as follows:

$$a_n \geq 0$$

$$t_n y(\mathbf{x}_n) - 1 + \xi_n \geq 0$$

$$a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} = 0$$

$$\mu_n \geq 0$$

$$\xi_n \geq 0$$

$$\mu_n \xi_n = 0$$

Quadratic Programming

$$L(\mathbf{w}, b, \xi, \mathbf{a}, \boldsymbol{\mu}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n$$

□ Optimizing wrt \mathbf{w}, b, ξ_n gives:

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \rightarrow \mathbf{w} = \sum_{n=1}^N a_n t_n \boldsymbol{\phi}(\mathbf{x}_n)$$

$$\frac{\partial L}{\partial b} = 0 \rightarrow \sum_{n=1}^N a_n t_n = 0$$

$$\frac{\partial L}{\partial \xi_n} = 0 \rightarrow a_n = C - \mu_n$$

□ Use these results to eliminate \mathbf{w}, b, ξ_n from $L(\mathbf{w}, b, \xi, \mathbf{a}, \boldsymbol{\mu})$.

Quadratic Programming

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

- This is identical to the separable case but with different constraints.
- $a_n \geq 0$ is required as they are Lagrange multipliers.
- $a_n = C - \mu_n$ together with $\mu_n \geq 0$ implies $a_n \leq C$.
- We thus need to maximize $\tilde{L}(\mathbf{a})$ wrt $a_n, n = 1, \dots, N$ subject to

$$\begin{aligned} 0 &\leq a_n \leq C \\ \sum_{n=1}^N a_n t_n &= 0 \end{aligned}$$

- Predictions are again made using: $y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$

Interpretation of the Solution

- As before, a subset of the data points may have $a_n = 0$, in which case they do not contribute to the predictive model $y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$
- The remaining data points constitute the support vectors. They are such that $a_n > 0$ and hence from $a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} = 0$ must satisfy:

$$t_n y(\mathbf{x}_n) = 1 - \xi_n$$

- If $a_n < C$ then $a_n = C - \mu_n$ implies $\mu_n > 0$ which from $\mu_n \xi_n = 0$ requires $\xi_n = 0$ and hence such points lie in the margin.
- Points with $a_n = C$ can lie inside the margin and can either be correctly classified if $\xi_n \leq 1$ or misclassified if $\xi_n > 1$.

Determining b

- To determine b in $y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$, we note that those support vectors for which $0 < a_n < C$ have $\xi_n = 0$ so that $t_n y(\mathbf{x}_n) = 1$ and hence will satisfy:

$$t_n \left(\sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) + b \right) = 1$$

- A numerically stable solution is obtained

$$b = \frac{1}{N_{\mathcal{M}}} \sum_{n \in \mathcal{M}} \left(t_n - \sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) \right)$$

- Here \mathcal{M} the sets of indices of data points having $0 < a_n < C$.

ν -SVM

- An alternative equivalent formulation of the SVM involves maximizing:

$$\tilde{L}(\mathbf{a}) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

subject to the following constraints:

$$0 \leq a_n \leq 1/N$$

$$\sum_{n=1}^N a_n t_n = 0$$

$$\sum_{n=1}^N a_n \geq \nu$$

- Scholkopf, B., A. Smola, R. C. Williamson, and P. L. Bartlett (2000). [New support vector algorithms](#). *Neural Computation* **12**(5), 1207–1245

ν -SVM

$$\tilde{L}(\mathbf{a}) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

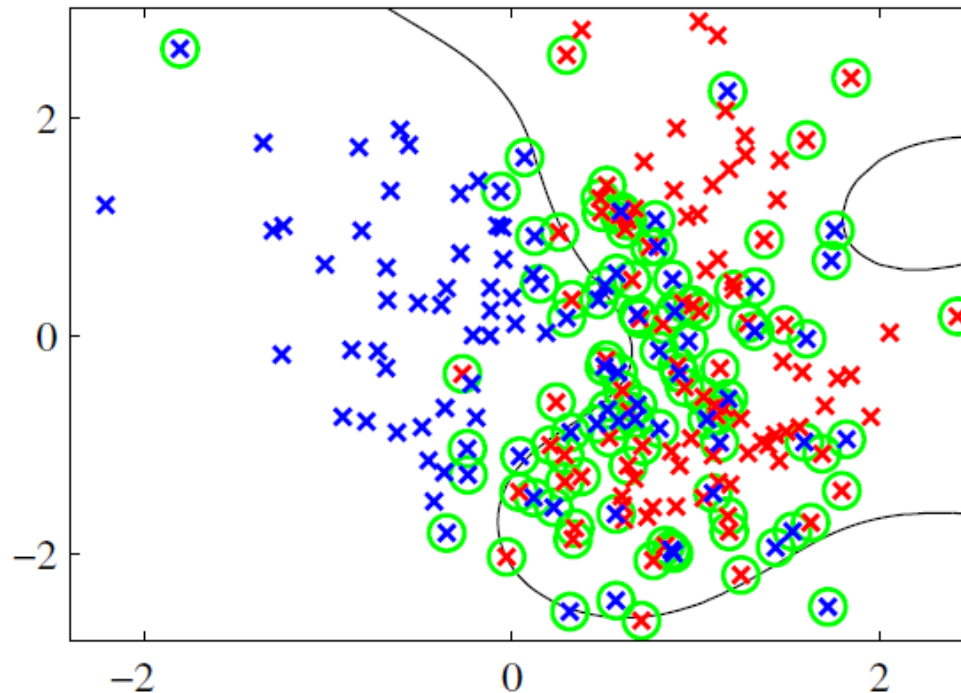
$$0 \leq a_n \leq 1/N, \quad \sum_{n=1}^N a_n t_n = 0, \quad \sum_{n=1}^N a_n \geq \nu$$

□ ν , which replaces C , can be interpreted as both

- ✓ an upper bound on the fraction of *margin errors* (points for which $\xi_n > 0$ and hence which lie on the wrong side of the margin boundary and which may or may not be misclassified) and
- ✓ a lower bound on the fraction of support vectors.

ν -SVM

- An example of the ν -SVM applied to a synthetic data set is shown below.
- Gaussian kernels $\exp(-\gamma||\mathbf{x} - \mathbf{x}'||^2)$ are used with $\gamma = 0.45$.



ν -SVM

- ❑ Although predictions for new inputs are made using only the support vectors, the training phase (i.e., the determination of the parameters \mathbf{a} and b) makes use of the whole data set.
- ❑ Thus it is important to have efficient algorithms for the quadratic programming problem.
- ❑ $\tilde{L}(\mathbf{a})$ is quadratic and so any local optimum will also be a global optimum provided the constraints define a convex region (which they do as a consequence of being linear).
- ❑ Direct solution of the quadratic programming problem using traditional techniques is often infeasible due to computation and memory requirements.

Chunking

- ❑ *Chunking* exploits the fact that the value of the Lagrangian is unchanged if we remove the rows and columns of the kernel matrix corresponding to Lagrange multipliers that have zero value.
 - ❑ This allows the full quadratic programming problem to be broken down into a series of smaller ones, whose goal is eventually to identify all of the nonzero Lagrange multipliers and discard the others.
 - ❑ Chunking can be implemented using *protected conjugate gradients*.
 - ❑ Although chunking reduces the size of the matrix in the quadratic function from the number of data points squared to approximately the number of nonzero Lagrange multipliers squared, even this may be too big to fit in memory for large-scale applications.
-
- Burges, C. J. C. (1998). [A tutorial on support vector machines for pattern recognition](#). *Knowledge Discovery and Data Mining* 2(2), 121–167.
 - Vapnik, V. N. (1982). [Estimation of dependences based on empirical data](#). Springer.

Decomposition Methods

- ❑ *Decomposition methods* also solve a series of smaller quadratic programming problems (each of these is of a fixed size), and so the technique can be applied to large data sets.
 - ❑ One of the most popular approaches to training support vector machines is called *sequential minimal optimization*.
 - ❑ It takes the concept of chunking to the extreme limit and considers just *two Lagrange multipliers at a time*. In this case, *the subproblem can be solved analytically*, avoiding quadratic programming.
 - ❑ Heuristics are given for choosing the pair of Lagrange multipliers in each step. The scaling with the number of data points is between linear and quadratic depending on the application.
-
- Osuna, E., R. Freund, and F. Girosi (1996). [Support vector machines: training and applications](#). A.I. Memo AIM-1602, MIT.
 - Platt, J. C. (1999). [Fast training of support vector machines using sequential minimal optimization](#). In B. Scholkopf, C. J. C. Burges, and A. J. Smola (Eds.), *Advances in Kernel Methods – Support Vector Learning*, pp. 185–208. MIT Press.

SVM and Curse of Dimensionality

- ❑ Kernel functions correspond to inner products in feature spaces that can have high, or even infinite, dimensionality.
- ❑ By working directly in terms of the kernel function, without introducing the feature space explicitly, does not mean that SVMs avoid the curse of dimensionality.
- ❑ There are constraints amongst the feature values **that restrict the effective dimensionality of feature space**. To see this consider a 2nd order polynomial kernel that we can expand in terms of its components

$$\begin{aligned}k(\mathbf{x}, \mathbf{z}) &= (1 + \mathbf{x}^T \mathbf{z})^2 = (1 + x_1 z_1 + x_2 z_2)^2 \\&= 1 + 2x_1 z_1 + 2x_2 z_2 + x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\&= (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1 x_2, x_2^2)(1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, \sqrt{2}z_1 z_2, z_2^2)^T \\&= \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{z})\end{aligned}$$

SVM and Curse of Dimensionality

- ❑ This kernel function therefore represents an inner product in a feature space having six dimensions, in which the mapping from input space to feature space is described by the vector function $\phi(x)$.
- ❑ However, the coefficients weighting these different features are constrained to have specific forms.
- ❑ Thus any set of points in the original 2D space x would be constrained to lie exactly on a 2D nonlinear manifold embedded in the six-dimensional feature space.

Probabilistic Predictions with SVM

- ❑ The SVM does not provide probabilistic outputs but makes classification decisions for new input vectors.
- ❑ Modifications to the SVM were introduced to allow the trade-off between false positive and false negative errors to be controlled.
- ❑ To address probabilistic predictions of the class label t for new inputs \mathbf{x} , Platt (2000) has proposed fitting a logistic sigmoid to the outputs of a previously trained SVM. The required conditional probability is

$$p(t = 1|\mathbf{x}) = \sigma (A y(\mathbf{x}) + B)$$

where $y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$.

- Veropoulos, K., C. Campbell, and N. Cristianini (1999). [Controlling the sensitivity of support vector machines](#). In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI99), Workshop ML3*, pp. 55–60.
- Platt, J. C., N. Cristianini, and J. Shawe-Taylor (2000). [Large margin DAGs for multiclass classification](#). In S. A. Solla, T. K. Leen, and K. R. Muller (Eds.), *Advances in Neural Information Processing Systems*, Volume 12, pp. 547–553. MIT Press.
- Tipping, M. (2001). [Sparse Bayesian learning and the relevance vector machine](#). *J. of Machine Learning Research* 1, 211–244.

Probabilistic Predictions with SVM

$$p(t = 1|\mathbf{x}) = \sigma (A y(\mathbf{x}) + B)$$

where $y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$.

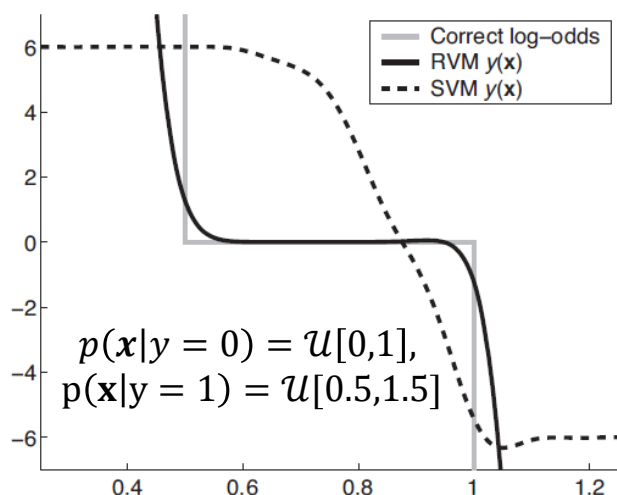
- ❑ Values for A and B are found by minimizing the cross-entropy error function defined by a training set consisting of pairs of $y(\mathbf{x}_n)$ and t_n .
- ❑ The data used to fit the sigmoid needs to be independent of that used to train the original SVM in order to avoid severe over-fitting. This two-stage approach is equivalent to assuming that the output $y(\mathbf{x})$ of the SVM represents the log-odds of \mathbf{x} belonging to class $t = 1$.
- ❑ The SVM can give a poor approximation to the posterior probabilities.
- Tipping, M. (2001). [Sparse Bayesian learning and the relevance vector machine](#). *J. of Machine Learning Research* 1, 211–244.

Probabilistic Response

- ❑ An SVM classifier produces a hard-labeling, $\hat{y}(x) = \text{sign}(f(x))$.
- ❑ A heuristic approach to producing confidence in our prediction is to interpret $f(x)$ as the log-odds ratio, $\log \frac{p(y = 1|x)}{p(y = 0|x)}$ and convert the output of an SVM to a probability using

$$p(y = 1|x, \theta) = \sigma(af(x) + b)$$

where a, b can be estimated by maximum likelihood on a separate validation set (using the training set to estimate a and b leads to severe overfitting.)



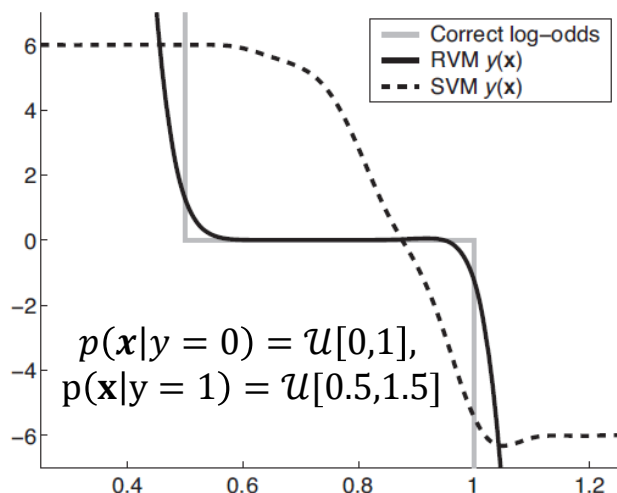
- The log-odds of class 1 over class 0 should be zero in $[0.5, 1.0]$, and infinite outside this region.
- We sampled 1000 points from the model, and then fit an RVM and an SVM with a Gaussian kernel of width 0.1.
- Both models can perfectly capture the decision boundary, and achieve a generalization error of 25%, which is Bayes optimal in this problem.
- The probabilistic output from the RVM is a good approximation to the true log-odds, but this is not the case for the SVM.

Probabilistic Response

- An SVM classifier produces a hard-labeling, $\hat{y}(\mathbf{x}) = \text{sign}(f(\mathbf{x}))$.
- A heuristic approach to producing confidence in our prediction is to interpret $f(\mathbf{x})$ as the log-odds ratio, $\log \frac{p(y = 1|\mathbf{x})}{p(y = 0|\mathbf{x})}$ and convert the output of an SVM to a probability using

$$p(y = 1|\mathbf{x}, \boldsymbol{\theta}) = \sigma(af(\mathbf{x}) + b)$$

where a, b can be estimated by maximum likelihood on a separate validation set (using the training set to estimate a and b leads to severe overfitting.)



- Tipping, M. (2001). [Sparse Bayesian learning and the relevance vector machine](#). *J. of Machine Learning Research* 1, 211–244.

Relation to Logistic Regression

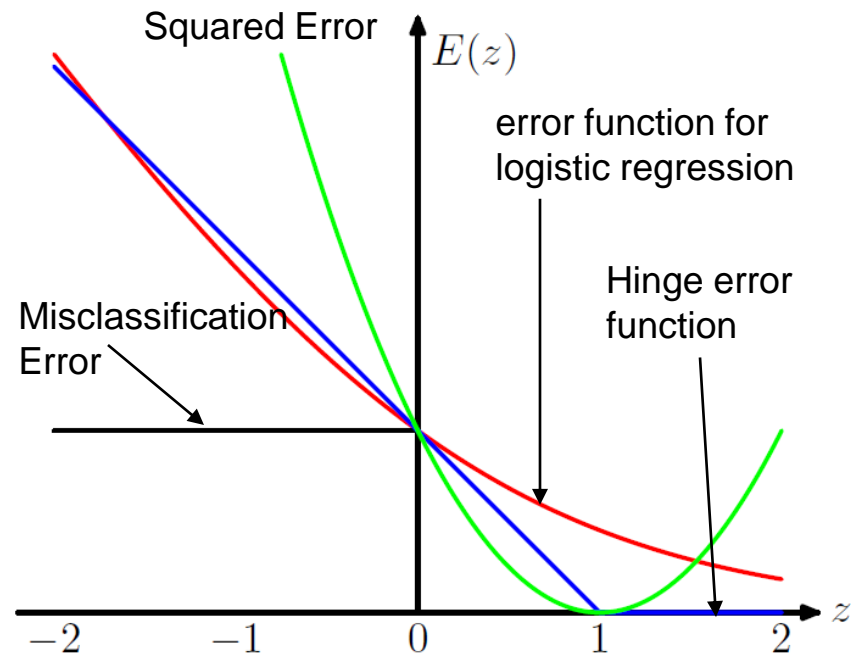
- We can re-cast the SVM for nonseparable distributions in terms of the minimization of a regularized error function.
- This will also allow us to highlight similarities/differences from the logistic regression.
- As a reminder for data points that are on the correct side of the margin boundary, and which therefore satisfy $y_n t_n \geq 1$, we have $\xi_n = 0$, and for the remaining points we have $\xi_n = 1 - y_n t_n$.
- Thus the objective function $C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$ can be written (up to a multiplicative constant) in the form

$$\sum_{n=1}^N E_{SV}(y_n t_n) + \lambda \|\mathbf{w}\|^2$$

where $\lambda = (2C)^{-1}$, and $E_{SV}(\cdot)$ is the *hinge error function* defined by $E_{SV}(y_n t_n) = [1 - y_n t_n]_+$ where $[\cdot]_+$ denotes the positive part.

Relation to Logistic Regression

- Plot of the 'hinge' error function in SVM, along with the error function for logistic regression, rescaled by a factor of $1/\ln(2)$ so that it passes through the point $(0, 1)$.
- When we considered the logistic regression model we found it convenient to work with target variable $t \in \{0, 1\}$.



For comparison with the SVM, we first reformulate maximum likelihood logistic regression using the target variable $t \in \{-1, 1\}$.

- To do this, we note that $p(t = 1|y) = \sigma(y)$ where $y(x) = \mathbf{w}^T \boldsymbol{\phi}(x) + b$ and $\sigma(y)$ is the logistic sigmoid function.

Error Function for Logistic Regression

- It follows that $p(t = -1|y) = 1 - \sigma(y) = \sigma(-y)$, where we have used the properties of the logistic sigmoid function, and so we can write $p(t|y) = \sigma(yt)$.
- We can construct an error function by taking the negative logarithm of the likelihood function that with a quadratic regularizer takes the form

$$\sum_{n=1}^N E_{LR}(y_n t_n) + \lambda \|\mathbf{w}\|^2, E_{LR}(y_n t_n) = \ln(1 + \exp(-yt))$$

- For comparison with other error functions, we can divide by $\ln(2)$ so that the error function passes through the point $(0, 1)$. This rescaled error function is also plotted in the earlier Fig. and has a similar form to the support vector error function.
- The key difference is that the flat region in $E_{SV}(y_n t_n)$ leads to sparse solutions. Both the logistic error and the hinge loss can be viewed as continuous approximations to the misclassification error.

Hinge Loss

- We have seen that the negative log-likelihood for the logistic regression model is

$$L_{nll}(y, t) = -\log p(t|\mathbf{x}, \mathbf{w}) = \log(1 + e^{-yt})$$

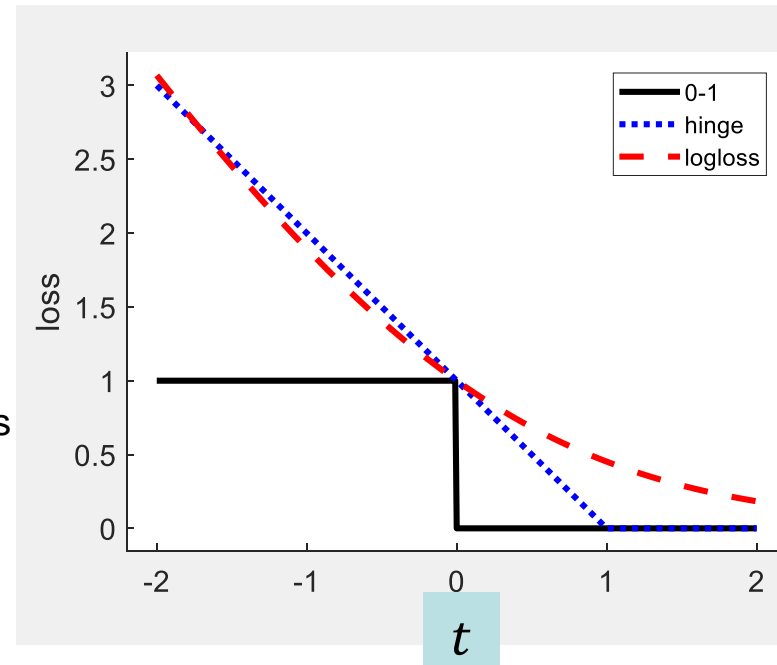
- It is a convex upper bound on the 0 – 1 risk of a binary classifier, where $y = f(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$ is the log odds ratio, and we have assumed the labels are $t \in \{1, -1\}$ rather than $\{0, 1\}$.

- We can replace the NLL loss with the **hinge loss**

$$L_{hinge}(y, t) = \max(0, 1 - yt) = (1 - yt)_+$$

The horizontal axis is the margin yt , the vertical axis is the loss. The log loss uses log base 2.

[hingeLossPlot.m](#)
from [PMTK3](#)



Hinge Loss

- Our optimization problem is as follows:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N (1 - t_i f(\mathbf{x}_i))_+$$

- Introducing slack variables we can write this as a quadratic program in $N + D + 1$ variables:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i, \quad \xi_i \geq 0, t_i(\boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w} + b) \geq 1 - \xi_i, i = 1, \dots, N$$

- We can eliminate \mathbf{w} , b and ξ_i , and just solve the N dual variables, which correspond to the Lagrange multipliers for the constraints. Standard solvers take $\mathcal{O}(N^3)$ time. **Sequential minimal optimization** takes $\mathcal{O}(N^2)$. Linear SVMs, take $\mathcal{O}(N)$ time to train.

- Platt, J. (1998). [Using analytic QP and sparseness to speed training of support vector machines](#). In *NIPS*.
- Joachims, T. (2006). [Training Linear SVMs in Linear Time](#). In *Proc. Of the Int'l Conf. on Knowledge Discovery and Data Mining*.
- Bottou, L., O. Chapelle, D. DeCoste, and J. Weston (Eds.) (2007). [Large Scale Kernel Machines](#). MIT Press.

SVMs for Classification

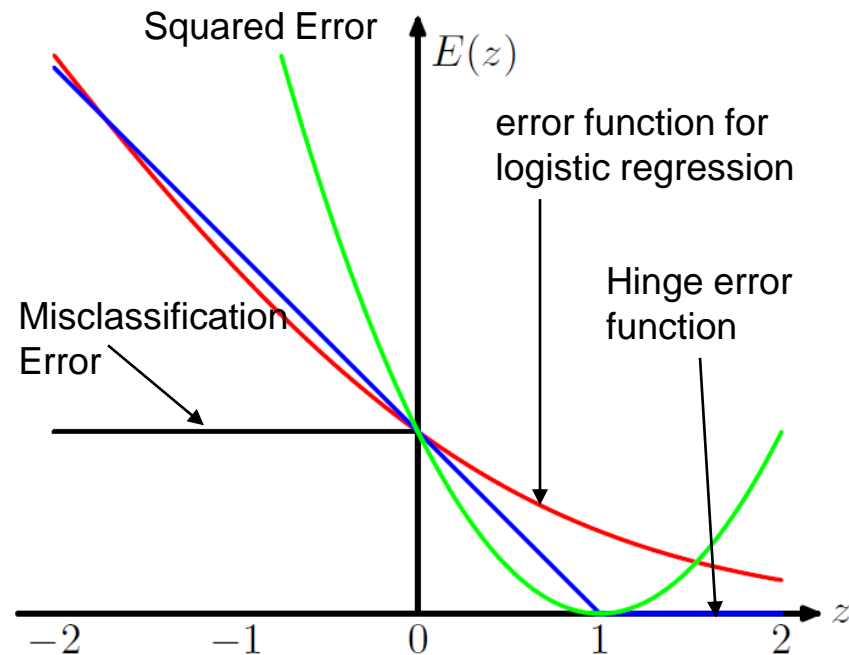
- We have shown that: $\hat{\mathbf{w}} = \sum_{i=1}^N \alpha_i t_i \boldsymbol{\phi}(\mathbf{x}_i)$, α is sparse
- The \mathbf{x}_i 's for which $\alpha_i > 0$ are called support vectors – these are points which are incorrectly classified or are classified correctly but are on or inside the margin.
- Making predictions:

$$\hat{y}(\mathbf{x}) = \text{sgn}(b + \hat{\mathbf{w}}^T \boldsymbol{\phi}(\mathbf{x})) = \text{sgn}(b + \sum_{i=1}^N \alpha_i k(\mathbf{x}_i, \mathbf{x}))$$

- It takes take $\mathcal{O}(sD)$ time to compute where $s \leq N$ is the number of support vectors.

Error Functions

- ❑ Another continuous error function that has sometimes been used to solve classification problems is the squared error.
- ❑ It has the property, however, of placing increasing emphasis on data points that are correctly classified but that are a long way from the decision boundary on the correct side.
- ❑ Such points will be strongly weighted at the expense of misclassified points, and so if the objective is to minimize the misclassification rate, then a monotonically decreasing error function would be a better choice.



Probabilistic Interpretation of SVM

- Thus the minimization objective for the Hinge loss function

$$C \sum_{i=1}^N (1 - t_i [b + \boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}])^+ + \frac{1}{2} \|\mathbf{w}\|^2$$

is the same as that based on logistic regression.

- The corresponding regularized loss takes the form:

$$\frac{1}{N} \sum_{i=1}^N (1 - t_i [b + \boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}])^+ + \frac{1}{NC} \|\mathbf{w}\|^2,$$

- Here the regularization parameter is $\lambda = \frac{1}{NC}$

- The first term above corresponds to the negative log likelihood:

$$\frac{1}{N} \sum_{i=1}^N -\log g(t_i [b + \boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}]) \quad , g(z) = \frac{1}{1 + e^{-z}} = \text{logistic function}$$

Probabilistic Interpretation of SVM

- The corresponding regularized loss takes the form:

$$\frac{1}{N} \sum_{i=1}^N (1 - t_i [b + \boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}])^+ + \frac{1}{NC} \|\mathbf{w}\|^2,$$

- One in the pursue of a probabilistic interpretation of SVM can show:

$$\begin{aligned} & \exp(-2(1 - t_i [b + \boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}])^+) \\ &= \int_0^\infty \frac{1}{\sqrt{2\pi\lambda_i}} \exp\left(-\frac{1}{2} \frac{(1 + \lambda_i - t_i [w_0 + \boldsymbol{\phi}(\mathbf{x}_i)^T \mathbf{w}])^2}{\lambda_i}\right) d\lambda_i \end{aligned}$$

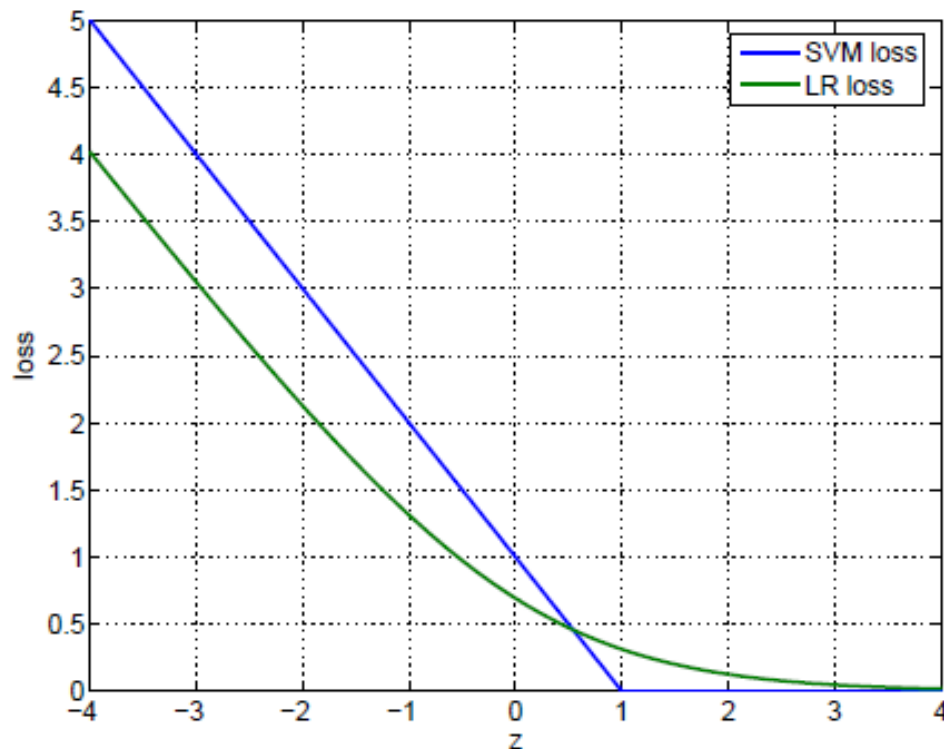
- This allows to fit the SVM with EM using λ_i as latent variables.

- Sollich, P. (2002). [Bayesian methods for support vector machines: evidence and predictive class probabilities](#). *Machine Learning* 46, 21– 52.
- Polson, N. and S. Scott (2011). [Data augmentation for support vector machines](#). *Bayesian Analysis* 6(1), 1–124.
- Franc, V., A. Zien, and B. Schoelkopf (2011). [Support vector machines as probabilistic models](#). In *Intl. Conf. on Machine Learning*.

Probabilistic Interpretation of SVM

□ Common objective the minimization of:

$$\frac{1}{N} \sum_{i=1}^N \text{Loss}(t_i[b + \phi(x_i)^T \mathbf{w}]) + \frac{1}{NC} \|\mathbf{w}\|^2,$$

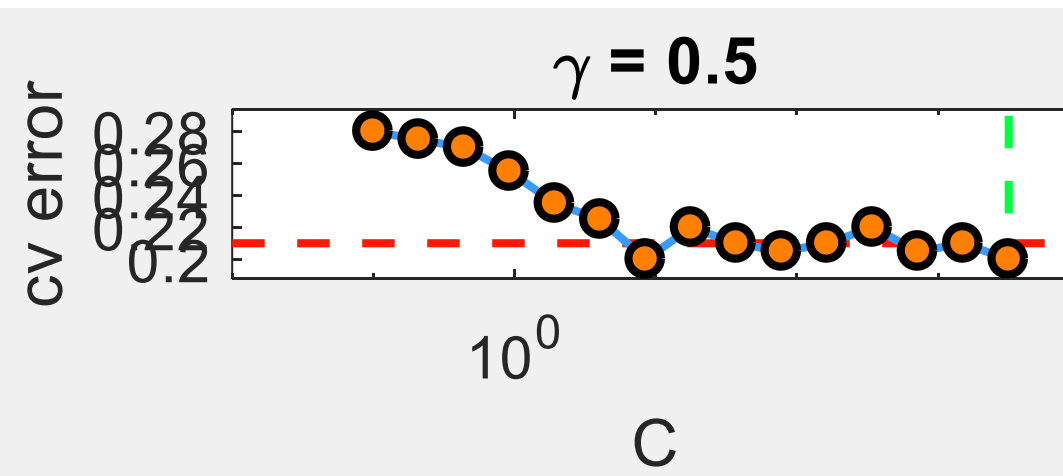
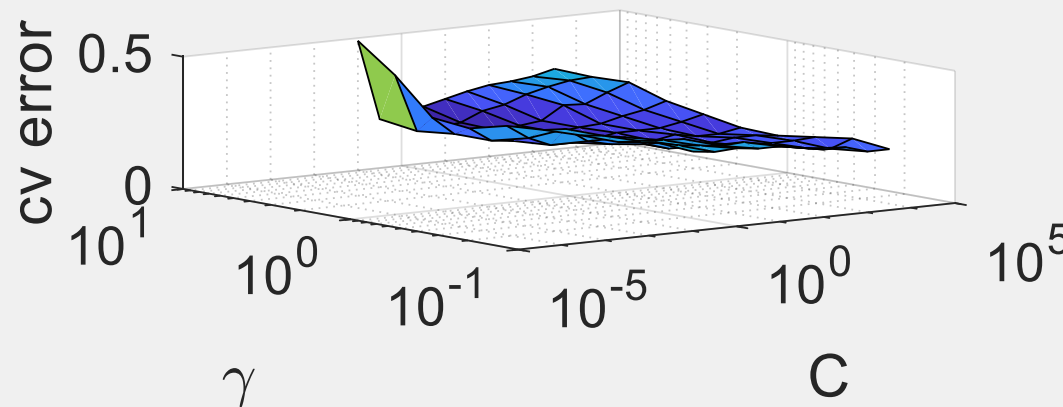


$$\text{SVM Loss: } \text{Loss}(z) = (1 - z)^+$$

$$\text{LR Loss: } \text{Loss}(z) = \frac{1}{1 + e^{-z}}$$

Choosing C

- Need to specify the kernel parameters and C . Typically C is chosen by cross-validation.
- C interacts quite strongly with the kernel parameters.
- Consider an RBF kernel with precision $\gamma = \frac{1}{(2\sigma)^2}$
- If $\gamma = 5$, corresponding to narrow kernels, we need heavy regularization, and hence small C (so $\lambda = 1/C$ is big).
- If $\gamma = 1$, a larger value of C should be used.



CV estimate of the 0-1 risk as a function of C and γ .

$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

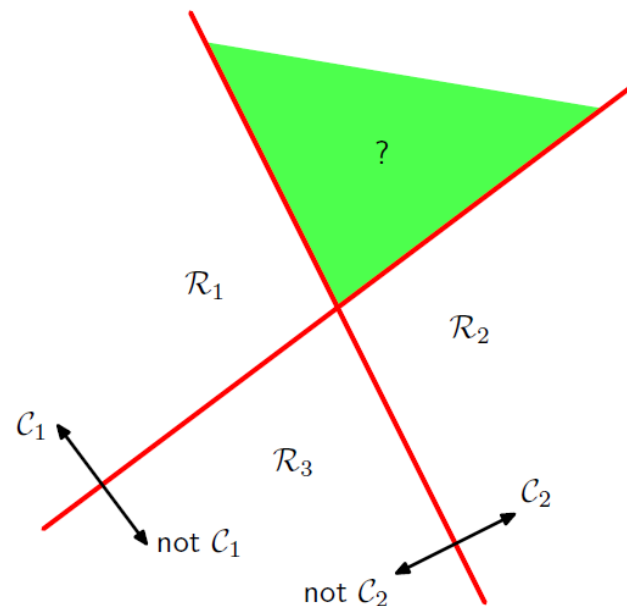
svmCGammaDemo.m
from PMTK3

Choosing C

- ❑ libsvm recommends using CV over a 2d grid with values $C \in \{2^{-5}, 2^{-3}, \dots, 2^{15}\}$ and $\gamma \in \{2^{-15}, 2^{-13}, \dots, 2^3\}$.
- ❑ Standardize the data first, for a spherical Gaussian kernel to make sense.
- ❑ To choose C efficiently, one can develop a path following the spirit of lars.
 - The basic idea is to start with λ large, so that the margin $1/||\mathbf{w}(\lambda)||$ is wide, and hence all points are inside of it and have $\alpha_i = 1$.
 - By slowly decreasing λ , a small set of points will move from inside the margin to outside, and their α_i values will change from 1 to 0, as they cease to be support vectors.
 - When λ is maximal, the function is completely smoothed, and no support vectors remain.
- Hastie, T., S. Rosset, R. Tibshirani, and J. Zhu (2004). [The entire regularization path for the support vector machine](#). *J. of Machine Learning Research* 5, 1391–1415.

Multiclass SVMs

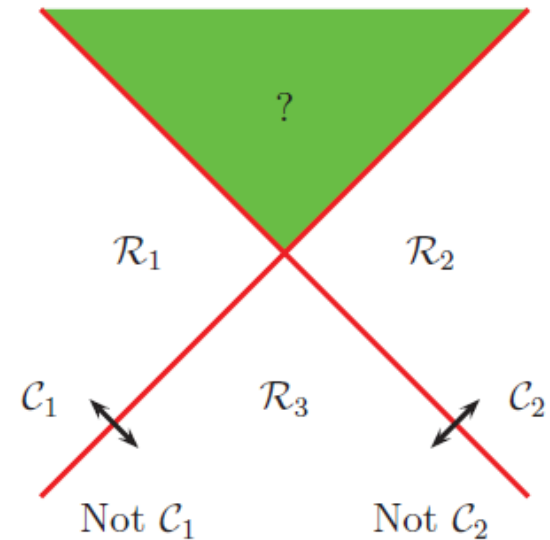
- ❑ The SVM is a two-class classifier. However, we often have to tackle $K > 2$ classes.
- ❑ Combine two-class SVMs in order to build a multiclass classifier.
- ❑ One approach (Vapnik, 1998) is to construct K separate SVMs, in which the k -th model $y_k(\mathbf{x})$ is trained using the data from class C_k as the positive examples and the data from the remaining $K - 1$ classes as the negative examples.
- ❑ This is known as the *one-versus-the-rest* approach.
- ❑ Using the decisions of the individual classifiers can lead to inconsistent results in which an input is assigned to multiple classes simultaneously.



- Vapnik, V. N. (1998). *Statistical learning theory*. Wiley.

Multiclass Classification

- ❑ The obvious approach is to use a **one-vs-all** approach. We train C binary classifiers, $f_c(x)$, where the data from class c is treated as positive, and the data from all the other classes as negative.
- ❑ This can result in ambiguously labeled regions.
- ❑ Can pick $\hat{y}(x) = \underset{c}{\operatorname{argmax}} f_c(x)$. This will not work since the f_c functions don't have comparable magnitudes.



The one-versus-rest approach.
The green region is predicted to be both class 1 and class 2.

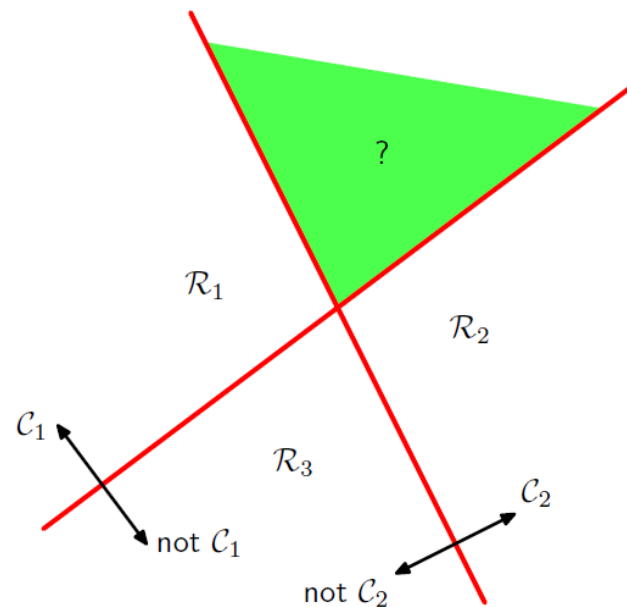
- Weston, J. and C. Watkins (1999). [Multi-class support vector machines](#). In *ESANN*.

Multiclass SVMs

- This problem is sometimes addressed by making predictions for new inputs \mathbf{x} using

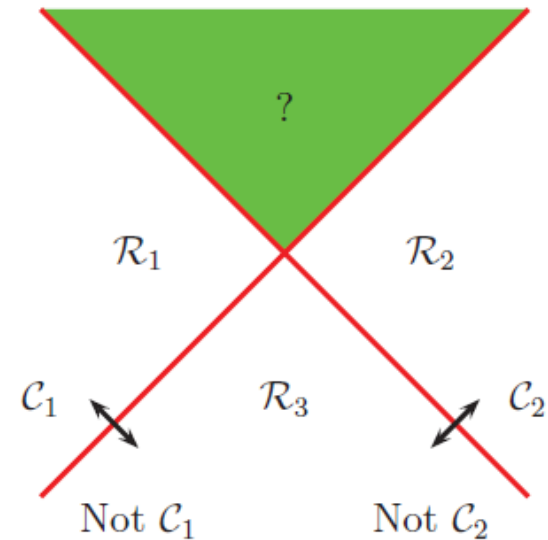
$$y(\mathbf{x}) = \max_k y_k(\mathbf{x})$$

- This heuristic approach suffers from the problem that the different classifiers were trained on different tasks, and **there is no guarantee that the $y_k(\mathbf{x})$ for different classifiers will have appropriate scales.**
- Another problem with the one-versus-the-rest approach is that **the training sets are imbalanced.**
- If e.g. we have 10 classes with equal numbers of training data points, then the individual classifiers are trained on data sets comprising 90% negative examples and 10% positive examples, and the symmetry of the original problem is lost.



Multiclass Classification

- ❑ Each binary subproblem is likely to suffer from the **class imbalance** problem.
- ❑ E.g. consider 10 equally represented classes. When training f_1 , we will have 10% positive examples and 90% negative examples, which can hurt performance.
- ❑ One can train all K classifiers simultaneously, but the resulting method takes $\mathcal{O}(K^2N^2)$ time, instead of $\mathcal{O}(KN^2)$.



The one-versus-rest approach.
The green region is predicted to be both class 1 and class 2.

Weston, J. and C. Watkins (1999). [Multi-class support vector machines](#). In *ESANN*.

Multiclass SVMs

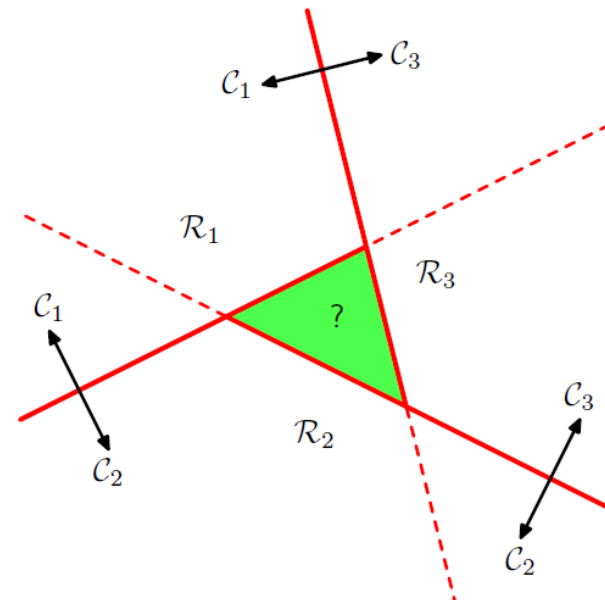
- A variant of the one-versus-the-rest scheme was proposed by Lee *et al.* (2001) who modify the target values so that the positive class has target $+1$ and the negative class has target $-1/(K - 1)$.
- Weston and Watkins (1999) define a single objective function for training all K SVMs simultaneously, based on maximizing the margin from each to remaining classes.

However, this can result in much slower training because, instead of solving K separate optimization problems each over N data points with an overall cost of $\mathcal{O}(KN^2)$, a single optimization problem of size $(K - 1)N$ must be solved giving an overall cost of $\mathcal{O}(K^2N^2)$

- Lee, Y., Y. Lin, and G. Wahba (2001). [Multicategory support vector machines](#). Technical Report 1040, Department of Statistics, University of Madison, Wisconsin.
- Weston, J. and C. Watkins (1999). [Multi-class support vector machines](#). In M. Verlysen (Ed.), *Proceedings ESANN'99, Brussels*. D-Facto Publications.

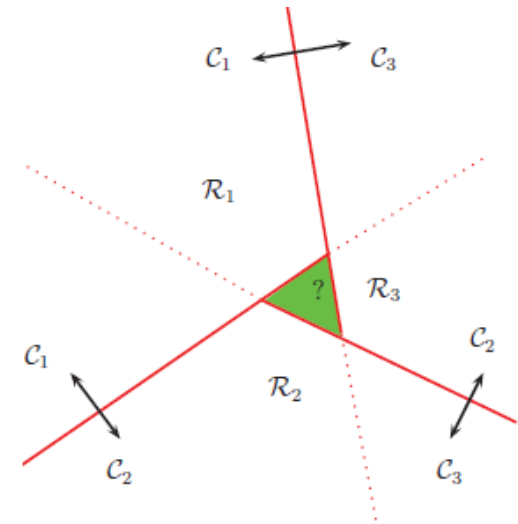
Multiclass SVMs

- ❑ One can train $K(K - 1)/2$ different 2-class SVMs on all possible pairs of classes, and then to classify test points according to which class has the highest number of ‘votes’, *one-versus-one*.
 - ❑ This can lead to ambiguities, for large K requires significantly more training time than the one-versus-the-rest approach and more computation is needed for predictions.
 - ❑ One can organize the pairwise classifiers into a directed acyclic graph leading to the *DAGSVM*. For K classes, the DAGSVM has a total of $K(K - 1)/2$ classifiers, and to classify a new test point only $K - 1$ pairwise classifiers need to be evaluated, with the particular classifiers used depending on which path through the graph is traversed.
- Platt, J. C., N. Cristianini, and J. Shawe-Taylor (2000). [Large margin DAGs for multiclass classification](#). In S. A. Solla, T. K. Leen, and K. R. Muller (Eds.), *Advances in Neural Information Processing Systems*, Volume 12, pp. 547–553. MIT Press.



Multiclass Classification

- ❑ Another approach is the **one-versus-one**.
- ❑ We train $K(K - 1)/2$ classifiers to discriminate all pairs $f_{c,c'}$. We then classify a point into the class which has the highest number of votes.
- ❑ This can also result in ambiguities.
- ❑ It takes $\mathcal{O}(K^2 N^2)$ time to train and $\mathcal{O}(K^2 N_{sv})$ to test each data point, where N_{sv} is the number of support vectors.
- ❑ Fundamental limitation of SVMs: they do not model uncertainty, so their output scores are not comparable across classes.



The one-versus-one approach.

- Allwein, E., R. Schapire, and Y. Singer (2000). [Reducing multiclass to binary: A unifying approach for margin classifiers](#). *J. of Machine Learning Research*, 113–141.

Multiclass SVMs

- ❑ A different approach to multiclass classification, based on [error-correcting output codes](#), was developed by Dietterich and Bakiri and applied to SVMs by Allwein *et al.*
 - ❑ This can be viewed as a generalization of the voting scheme of the one-versus-one approach in which [more general partitions of the classes are used to train the individual classifiers](#). The K classes themselves are represented as particular sets of responses from the two-class classifiers chosen, and together with a suitable decoding scheme, this gives robustness to errors and to ambiguity in the outputs of the individual classifiers.
 - ❑ The one-versus-the-rest approach is the most widely used in spite of its ad-hoc formulation and its practical limitations.
-
- Dietterich, T. G. and G. Bakiri (1995). [Solving multiclass learning problems via error-correcting output codes](#). *Journal of Artificial Intelligence Research* **2**, 263–286.
 - Scholkopf, B., J. Platt, J. Shawe-Taylor, Allwein, E. L., R. E. Schapire, and Y. Singer (2000). [Reducing multiclass to binary: a unifying approach for margin classifiers](#). *Journal of Machine Learning Research* **1**, 113–141.

SVMs for Unsupervised Learning

- ❑ There are also *single-class SVMs*, which solve an unsupervised learning problem related to probability density estimation.
 - ❑ Instead of modelling the density of data, these methods find a smooth boundary enclosing a region of high density. The boundary is chosen to represent a quantile of the density i.e. the probability that a data point drawn from the distribution will land inside that region is given by a number between 0 and 1 that is specified in advance.
 - ❑ This is a more restricted problem than estimating the full density. Two approaches to this problem using SVMs have been proposed.
-
- Tax, D. and R. Duin (1999). [Data domain description by support vectors](#). In M. Verleysen (Ed.), *Proceedings European Symposium on Artificial Neural Networks, ESANN*, pp. 251–256. D. Facto Press.
 - Scholkopf, B., J. Platt, J. Shawe-Taylor, A. Smola, and R. C. Williamson (2001). [Estimating the support of a high-dimensional distribution](#). *Neural Computation* **13**(7), 1433–1471.

Multiclass SVMs

- ❑ Scholkopf *et al.* (2001) tries to find a hyperplane that separates all but a fixed fraction ν of the training data from the origin while at the same time maximizing the distance (margin) of the hyperplane from the origin.
 - ❑ Tax and Duin (1999) look for the smallest sphere in feature space that contains all but a fraction ν of the data points.
 - ❑ For kernels $k(x, x')$ that are functions only of $x - x'$, the two algorithms are equivalent.
-
- Tax, D. and R. Duin (1999). [Data domain description by support vectors](#). In M. Verleysen (Ed.), *Proceedings European Symposium on Artificial Neural Networks, ESANN*, pp. 251–256. D. Facto Press.
 - Scholkopf, B., J. Platt, J. Shawe-Taylor, A. Smola, and R. C. Williamson (2001). [Estimating the support of a high-dimensional distribution](#). *Neural Computation* **13**(7), 1433–1471.

SVMs for Regression

- We now extend support vector machines to regression problems.
- In simple linear regression, we minimize a regularized error function given by

$$\frac{1}{2} \sum_{n=1}^N \{y_n - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- To obtain sparse solutions, the quadratic error function is replaced by an ϵ -insensitive error function (Vapnik, 1995), which gives zero error if the absolute difference between the prediction $y(\mathbf{x})$ and the target t is less than ϵ where $\epsilon > 0$.
- A simple example of an ϵ -insensitive error function, having a linear cost associated with errors outside the insensitive region, is given by

$$\mathbb{E}_{\epsilon}(y(\mathbf{x}) - t) = \begin{cases} 0, & \text{if } |y(\mathbf{x}) - t| < \epsilon: \\ |y(\mathbf{x}) - t| - \epsilon, & \text{otherwise} \end{cases}$$

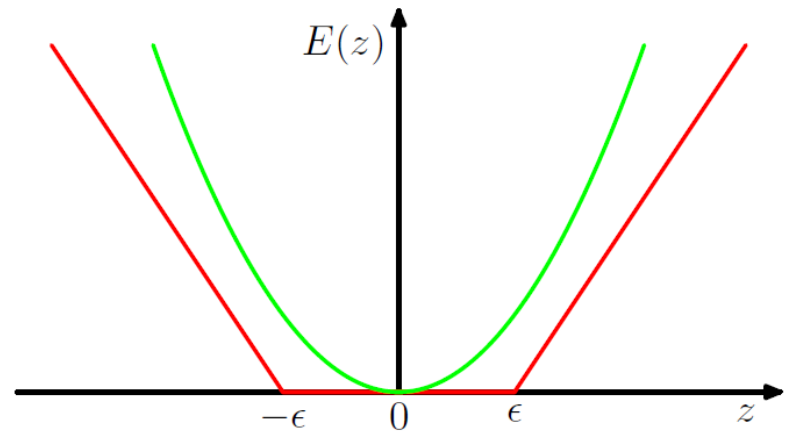
- Vapnik, V. N. (1995). *The nature of statistical learning theory*. Springer.
Machine Learning, University of Notre Dame, Notre Dame, IN, USA (Spring 2019, N. Zabarar)

SVMs for Regression

- We therefore minimize a regularized error function given by

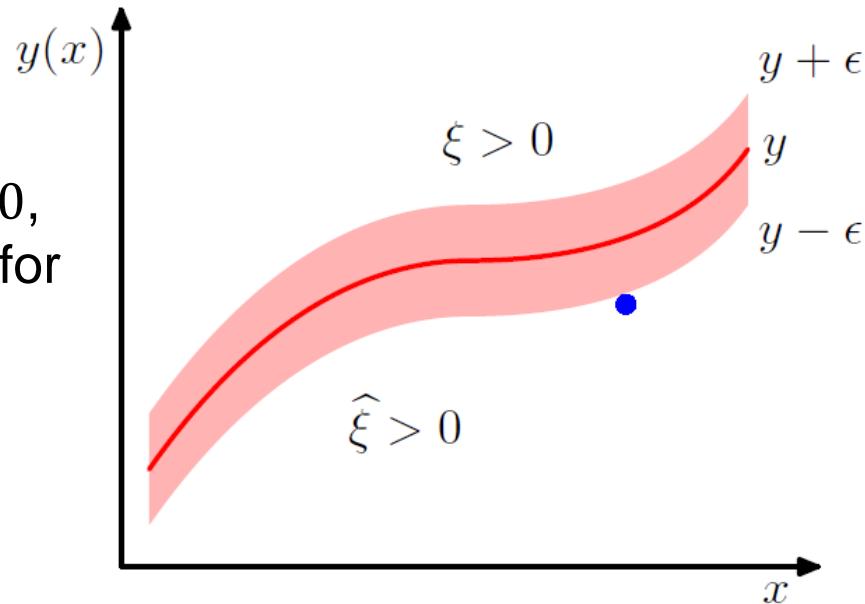
$$C \sum_{n=1}^N E_{\epsilon}(y(\mathbf{x}_n) - t_n) + \frac{1}{2} \|\mathbf{w}\|^2$$

where $y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$.



- The (inverse) regularization parameter, C , appears in front of the error term.

- For each data point \mathbf{x}_n , we now need two slack variables $\xi_n \geq 0$ and $\hat{\xi}_n \geq 0$, where $\xi_n \geq 0$ corresponds to a point for which $t_n > y(\mathbf{x}_n) + \epsilon$, and $\hat{\xi}_n > 0$, corresponds to a point for which $t_n < y(\mathbf{x}_n) - \epsilon$.



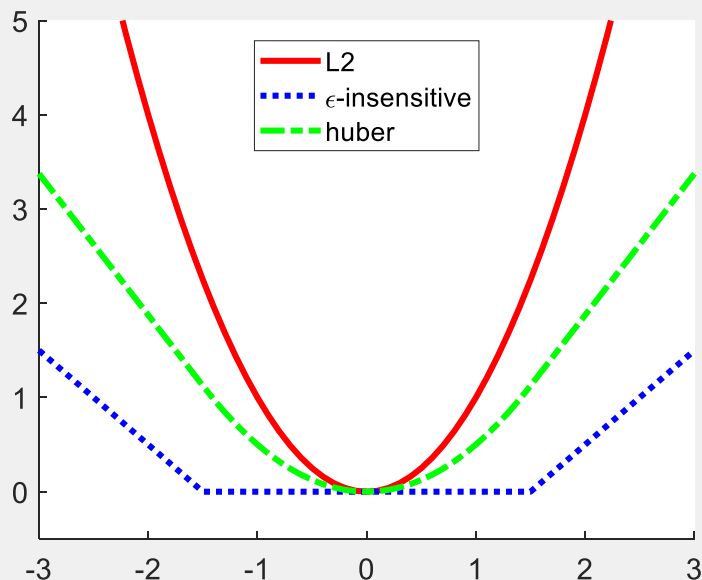
- Schoelkopf, B. and A. Smola (2002). [Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond](#). MIT Press.

SVMs for Regression

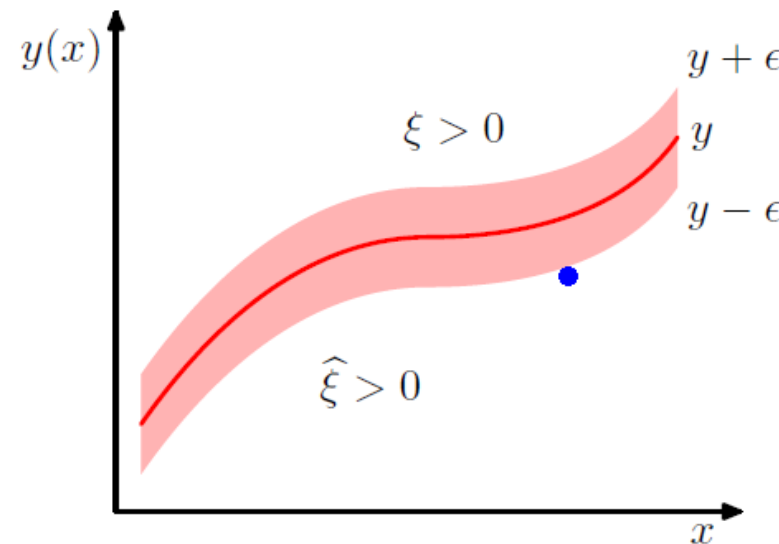
□ In our approach to regression, \mathbf{w} depends on all training inputs. We are here looking for sparse approximations.

□ The ϵ –insensitive loss function: $L_\epsilon(y, \hat{y}) \equiv \begin{cases} 0 & \text{if } |y - \hat{y}| < \epsilon \\ |y - \hat{y}| - \epsilon & \text{if } |y - \hat{y}| > \epsilon \end{cases}$
leading to the following objective function:

$$J = C \sum_{i=1}^N L_\epsilon(y_i, \hat{y}_i) + \frac{1}{2} \|\mathbf{w}\|^2, \hat{y}_i = f(\mathbf{x}_i) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + w_0, C = 1/\lambda$$



[huberLossDemo.m](#)
from [PMTK3](#)



SVMs for Regression

- The condition for a target point to lie inside the ε -tube is that $y_n - \varepsilon \leq t_n \leq y_n + \varepsilon$ where $y_n = y(\mathbf{x}_n)$. Introducing the slack variables allows points to lie outside the tube provided the slack variables are nonzero, and the corresponding conditions are

$$t_n \leq y(\mathbf{x}_n) + \varepsilon + \xi_n$$

$$t_n \geq y(\mathbf{x}_n) - \varepsilon - \hat{\xi}_n$$

- The error function for support vector regression can then be written as

$$C \sum_{n=1}^N (\hat{\xi}_n + \xi_n) + \frac{1}{2} \|\mathbf{w}\|^2$$

which must be minimized subject to the constraints above as well as the constraints for the Lagrange multipliers $a_n \geq 0$ and $\hat{a}_n \geq 0$, $\mu_n \geq 0$ and $\hat{\mu}_n \geq 0$.

$$\begin{aligned} L = & C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N (\mu_n \xi_n + \hat{\mu}_n \hat{\xi}_n) \\ & - \sum_{n=1}^N a_n (\varepsilon + \xi_n + y_n - t_n) - \sum_{n=1}^N \hat{a}_n (\varepsilon + \hat{\xi}_n - y_n + t_n) \end{aligned}$$

SVMs for Regression

$$L = C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N (\mu_n \xi_n + \hat{\mu}_n \hat{\xi}_n) - \sum_{n=1}^N a_n (\epsilon + \xi_n + y_n - t_n) - \sum_{n=1}^N \hat{a}_n (\epsilon + \hat{\xi}_n - y_n - t_n)$$

□ We substitute $y(\mathbf{x})$ and set the derivatives wrt $\mathbf{w}, b, \xi_n, \hat{\xi}_n$ equal to 0.

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \rightarrow \mathbf{w} = \sum_{n=1}^N (a_n - \hat{a}_n) \boldsymbol{\phi}(\mathbf{x}_n)$$

$$\frac{\partial L}{\partial b} = 0 \rightarrow \sum_{n=1}^N (a_n - \hat{a}_n) = 0$$

$$\frac{\partial L}{\partial \xi_n} = 0 \rightarrow a_n + \mu_n = C$$

$$\frac{\partial L}{\partial \hat{\xi}_n} = 0 \rightarrow \hat{a}_n + \hat{\mu}_n = C$$

SVM Regression

- Eliminating the corresponding variables, we arrive at the dual problem

$$\begin{aligned}\tilde{L}(\mathbf{a}, \hat{\mathbf{a}}) = & -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m)k(\mathbf{x}_n, \mathbf{x}_m) \\ & -\varepsilon \sum_{n=1}^N (a_n + \hat{a}_n) + \sum_{n=1}^N (a_n - \hat{a}_n)t_n\end{aligned}$$

- From $a_n + \mu_n = C$, $\hat{a}_n + \hat{\mu}_n = C$ and $\mu_n \geq 0$ and $\hat{\mu}_n \geq 0$, the following box constraints are applied

$$0 \leq a_n \leq C \text{ and } 0 \leq \hat{a}_n \leq C$$

in addition to $a_n \geq 0$ and $\hat{a}_n \geq 0$, $\mu_n \geq 0$ and $\hat{\mu}_n \geq 0$

- The predictions for new inputs are made using:

$$y(\mathbf{x}) = \sum_{n=1}^N (a_n - \hat{a}_n)k(\mathbf{x}, \mathbf{x}_n) + b$$

KKT Conditions

- The KKT conditions are:

$$a_n(\varepsilon + \xi_n + y_n - t_n) = 0$$

$$\hat{a}_n(\varepsilon + \hat{\xi}_n - y_n + t_n) = 0$$

$$(C - a_n)\xi_n = 0$$

$$(C - \hat{a}_n)\hat{\xi}_n = 0$$

- We note that a_n can only be nonzero if $\varepsilon + \xi_n + y_n - t_n = 0$, which implies that the data point either lies on the upper boundary of the ε -tube ($\xi_n = 0$) or lies above the upper boundary ($\xi_n = 0$).
- Similarly, a nonzero value for \hat{a}_n implies $\varepsilon + \hat{\xi}_n - y_n + t_n = 0$ and such points must lie either on or below the lower boundary of the ε -tube.

KKT Conditions

- The KKT conditions are:

$$a_n(\varepsilon + \xi_n + y_n - t_n) = 0$$

$$\hat{a}_n(\varepsilon + \hat{\xi}_n - y_n + t_n) = 0$$

$$(C - a_n)\xi_n = 0$$

$$(C - \hat{a}_n)\hat{\xi}_n = 0$$

- Furthermore, the two constraints $\varepsilon + \xi_n + y_n - t_n = 0$ and $\varepsilon + \hat{\xi}_n - y_n + t_n = 0$ are incompatible, as is easily seen by adding them together and noting that ξ_n and $\hat{\xi}_n$ are nonnegative while ε is strictly positive, and so for every data point x_n , either a_n or \hat{a}_n (or both) must be zero.
- The support vectors are those data points that contribute to predictions i.e. those for which either $a_n \neq 0$ or $\hat{a}_n \neq 0$. These are points that lie on the boundary of the -tube or outside the tube.

KKT Conditions

- All points within the tube have $a_n = \hat{a}_n = 0$, i.e. we obtain again a sparse solution.
- The parameter b can be found by considering a data point for which $0 < a_n < C$, which from $(C - a_n)\xi_n = 0$ must have $\xi_n = 0$, and from $a_n(\varepsilon + \xi_n + y_n - t_n) = 0$ must therefore satisfy $\varepsilon + y_n - t_n = 0$. Using $y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$ and $\mathbf{w} = \sum_{n=1}^N (a_n - \hat{a}_n) \boldsymbol{\phi}(\mathbf{x}_n)$ and solving for b , we obtain

$$b = t_n - \varepsilon - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) = t_n - \varepsilon - \sum_{m=1}^N (a_m - \hat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m)$$

KKT Conditions

- There is an alternative formulation of the SVM for regression. Instead of fixing the width ε of the insensitive region, we fix instead a parameter ν that bounds the fraction of points lying outside the tube. We maximize:

$$\tilde{L}(\mathbf{a}, \hat{\mathbf{a}}) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m) + \sum_{n=1}^N (a_n - \hat{a}_n) t_n$$

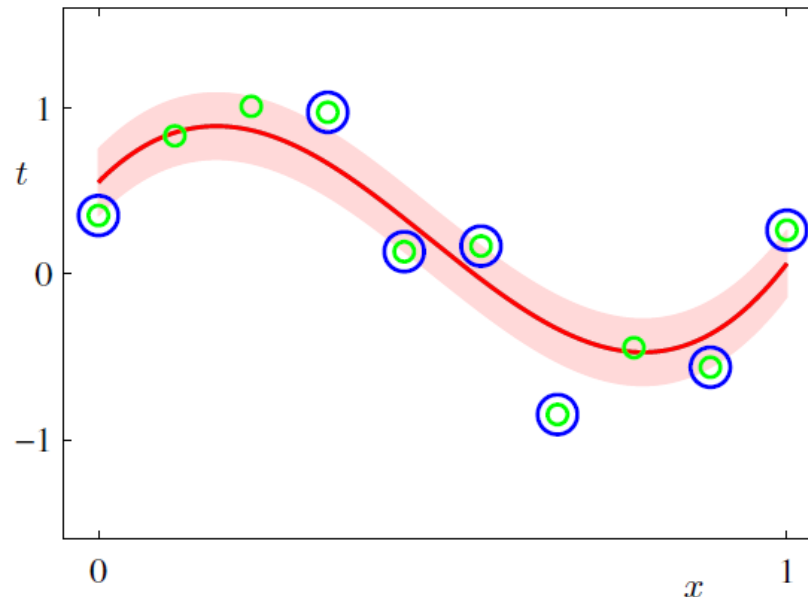
subject to:

$$0 \leq a_n \leq \frac{C}{N}, 0 \leq \hat{a}_n \leq \frac{C}{N}, \sum_{n=1}^N (a_n - \hat{a}_n) = 0, \sum_{n=1}^N (a_n + \hat{a}_n) \leq \nu C$$

- It can be shown that there are at most νN data points falling outside the insensitive tube, while at least νN data points are support vectors and so lie either on the tube or outside it.

ν -SVM for Regression

- The use of a support vector machine to solve a regression problem is illustrated below. Here the parameters ν and C have been chosen by hand. Their values would typically be determined by cross validation.



- Illustration of the ν -SVM for regression applied to the sinusoidal synthetic data set using Gaussian kernels. The predicted regression curve is shown by the red line, and the ϵ -insensitive tube corresponds to the shaded region. Also, the data points are shown in green, and those with support vectors are indicated by blue circles.

Computational Learning Theory

- ❑ SVMs have been motivated and analyzed using *computational learning theory*, also sometimes called *statistical learning theory*.
 - ❑ Valiant formulated the *probably approximately correct*, or PAC, learning framework.
 - ❑ The goal of PAC is to learn how large a data set needs to be in to give good generalization. It also gives bounds for the cost of learning.
 - ❑ Suppose a data set \mathcal{D} of size N is drawn from $p(x, t)$ where x is the input variable and t represents the class label, and restrict to ‘noise free’ situations in which t are determined by some (unknown) deterministic function $t = g(x)$.
- Valiant, L. G. (1984). [A theory of the learnable](#). *Communications of the Association for Computing Machinery* **27**, 1134–1142.
 - Anthony, M. and N. Biggs (1992). *An Introduction to Computational Learning Theory*. Cambridge University Press.
 - Kearns, M. J. and U. V. Vazirani (1994). *An Introduction to Computational Learning Theory*. MIT Press.
 - Vapnik, V. N. (1995). *The nature of statistical learning theory*. Springer.
 - Vapnik, V. N. (1998). *Statistical learning theory*. Wiley.
 - McAllester, D. A. (2003). [PAC-Bayesian stochastic model selection](#). *Machine Learning* **51**(1), 5–21.

Computational Learning Theory

- In PAC learning we say that a function $f(x; \mathcal{D})$, drawn from a space \mathcal{F} of such functions on the basis of the training set \mathcal{D} , has good generalization if its expected error rate is below some pre-specified threshold, so that

$$\mathbb{E}_{x,t}[\mathbb{I}(f(x; \mathcal{D}) = t)] < \varepsilon$$

- The quantity on the left-hand side is a random variable, because it depends on the training set \mathcal{D} , and the PAC framework requires that the Eq. above holds with probability greater than $1 - \delta$, for a data set \mathcal{D} drawn randomly from $p(x, t)$.
- Here δ is another pre-specified parameter, and the terminology ‘probably approximately correct’ comes from the requirement that with high probability (greater than $1 - \delta$), the error rate be small (less than ε).

Computational Learning Theory

$$\mathbb{E}_{\mathbf{x},t}[\mathbb{I}(\mathbf{f}(\mathbf{x};\mathcal{D}) = \mathbf{t})] < \varepsilon$$

- For a given choice of model space \mathcal{F} , and for given parameters ε and δ , PAC learning aims to provide bounds on the minimum size N of data set needed to meet this criterion.
- A key quantity in PAC learning is the *Vapnik-Chervonenkis dimension*, or VC dimension, which provides a measure of the complexity of a space of functions, and which allows the PAC framework to be extended to spaces containing an infinite number of functions.
- The bounds derived within the PAC framework are often described as worst case, because they apply to *any* choice for the distribution $p(\mathbf{x}, t)$, so long as both the training and the test examples are drawn (independently) from the same distribution, and for *any* choice for the function $\mathbf{f}(\mathbf{x})$ so long as it belongs to \mathcal{F} .

Computational Learning Theory

- ❑ In real-world applications of machine learning, we deal with distributions that have significant regularity, for example in which large regions of input space carry the same class label.
 - ❑ The PAC bounds are very conservative, in other words they strongly over-estimate the size of data sets required to achieve a given generalization performance. For this reason, PAC bounds have found few, if any, practical applications.
 - ❑ One attempt to improve the tightness of the PAC bounds is the *PAC-Bayesian framework*, which considers a distribution over the space \mathcal{F} of functions, somewhat analogous to the prior in a Bayesian treatment. This still considers any possible choice for $p(\mathbf{x}, \mathbf{t})$, and so although the bounds are tighter, they are still very conservative.
-
- McAllester, D. A. (2003). [PAC-Bayesian stochastic model selection](#). *Machine Learning* **51**(1), 5–21.

Relevance Vector Machines

Relevance Vector Machines

- ❑ Relevance Vector Machines (RVM) were introduced by [M. Tipping](#) in [Tipping, M. E. \(2001\). Sparse Bayesian learning and the relevance vector machine. Journal of Machine Learning Research 1, 211–244.](#)
- ❑ It is a *Bayesian regression* technique on a *generalized linear model* that produces *sparse representations* (in the sense that many of the coefficients of the generalized linear model turn out to be zero).
- ❑ We will first describe it on one output dimension but can generalize it to many dimensions.
- ❑ One can couple it with an adaptive input decomposition algorithm and apply it to uncertainty quantification tasks. [See I. Bilonis and N. Zabaras. Multidimensional adaptive relevance vector machines for uncertainty quantification. SIAM Journal for Scientific Computing, Vol. 34, No. 6, pp. B881–B908 2012](#) for further details.

Relevance Vector Machines

- ❑ RVM addresses many of the limitations of SVM:
 - SVM outputs are decisions rather than posterior probabilities.
 - Extension to $K > 2$ classes in SVM is problematic.
 - There are complexity parameters in SVM to be found by cross-validation.
 - Predictions in SVM are based on kernel functions centered at the training data that are required to be positive definite.

RVM: Likelihood

- We assume that we observe the data:

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)} = f(\mathbf{x}^{(i)})\}_{i=1}^n,$$

and that we wish to learn the function $f(\cdot)$.

- Pick a set of M basis functions $\boldsymbol{\phi}: \mathbf{R}^k \rightarrow \mathbf{R}^M$:

$$\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})).$$

- We will approximate $f(\cdot)$ by:

$$\hat{f}(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + \epsilon = \sum_{j=1}^M w_j \phi_j(\mathbf{x}) + \epsilon,$$

where ϵ is a random variable representing noise.

- For simplicity, we choose to work with Gaussian noise with zero mean and inverse variance β , so the **likelihood** is:

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}), \beta^{-1})$$

RVM: Prior

- The essence of RVM is the following **prior on the weights**:

$$p(\mathbf{w}|\boldsymbol{\alpha}) = \prod_{j=1}^M p(w_j|\alpha_j),$$

with $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_M)$ and:

$$p(w_j|\alpha_j) = \mathcal{N}(w_j|0, \alpha_j^{-1}).$$

- The characteristic of this prior is that:

$$\lim_{\alpha_j \rightarrow \infty} p(w_j|\alpha_j) = \delta(w_j).$$

- This means that the j -th basis function can be removed from the model if its α_j is very big. This is because, its weights will be very sharply picked about zero.
- To complete the model, we also need to assign a prior $p(\boldsymbol{\beta}, \boldsymbol{\alpha})$. We manage to get rid of it by employing the *evidence approximation*.

RVM: Posterior

□ The **posterior of the weights** is:

$$p(\mathbf{w}|\mathcal{D}, \beta, \alpha) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \Sigma),$$

where:

$$\Sigma = (\text{diag}(\alpha) + \beta \Phi^T \Phi)^{-1},$$
$$\mathbf{m} = \beta \Sigma \Phi^T \mathbf{t},$$

with $\mathbf{t} = (t^{(1)}, \dots, t^{(n)})$ is the *vector of observations*, $\Phi \in \mathbf{R}^{N \times M}$ is the *design matrix*:

$$\Phi_{ij} = \phi_j(\mathbf{x}^{(i)}),$$

and $\text{diag}(\alpha)$ is a diagonal matrix with α on the diagonal.

Proof:

First, notice that we can write:

$$p(\mathbf{w}|\alpha) = \prod_{j=1}^M p(w_j|\alpha_j) = \prod_{j=1}^M \mathcal{N}(w_j|0, \alpha_j^{-1}) = \mathcal{N}_M(\mathbf{w}|\text{diag}(\alpha)^{-1}).$$

RVM: Posterior

Proof (continuation):

Also, notice that the likelihood of the data (under an independence assumption) is given by:

$$p(\mathcal{D}|\mathbf{w}, \beta) = \prod_{i=1}^n \mathcal{N}(t^{(i)} | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(i)}), \beta^{-1}) = \mathcal{N}_N(\mathbf{y} | \boldsymbol{\Phi} \mathbf{w}, \beta^{-1} \mathbf{I}_N),$$

where \mathbf{I}_N is the N -dimensional unit matrix and we have used the fact that:

$$\begin{aligned} \sum_{i=1}^N \| \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(i)}) - t^{(i)} \|^2 &= \sum_{i=1}^N (\boldsymbol{\phi}(\mathbf{x}^{(i)}) \mathbf{w} - t^{(i)})^T (\boldsymbol{\phi}(\mathbf{x}^{(i)}) \mathbf{w} - t^{(i)}) \\ &= (\boldsymbol{\Phi} \mathbf{w} - \mathbf{y})^T (\boldsymbol{\Phi} \mathbf{w} - \mathbf{y}). \end{aligned}$$

Finally, we have:

$$p(\mathbf{w} | \mathcal{D}, \boldsymbol{\alpha}, \beta) \propto p(\mathcal{D} | \mathbf{w}, \beta) p(\mathbf{w} | \boldsymbol{\alpha}) = \mathcal{N}_N(\mathbf{y} | \boldsymbol{\Phi} \mathbf{w}, \beta^{-1} \mathbf{I}_N) \mathcal{N}_M(\mathbf{w} | \text{diag}(\boldsymbol{\alpha})^{-1}).$$

RVM: Posterior

Proof (continuation):

Remember that the posterior is:

$$p(\mathbf{w}|\mathcal{D}, \boldsymbol{\alpha}, \beta) \propto p(\mathcal{D}|\mathbf{w}, \beta)p(\mathbf{w}|\boldsymbol{\alpha}) = \mathcal{N}_N(\mathbf{y}|\boldsymbol{\Phi}\mathbf{w}, \beta^{-1}\mathbf{I}_N)\mathcal{N}_M(\mathbf{w}|\text{diag}(\boldsymbol{\alpha})^{-1}).$$

Now, we just have to look at it and complete the square...

You can simply drop this guy!

We have for the exponential inside the right hand side:

$$\begin{aligned} & (\boldsymbol{\Phi}\mathbf{w} - \mathbf{t})^T(\beta\mathbf{I}_N)(\boldsymbol{\Phi}\mathbf{w} - \mathbf{t}) + \mathbf{w}^T\text{diag}(\boldsymbol{\alpha})\mathbf{w} = \\ & \mathbf{w}^T(\beta\boldsymbol{\Phi}^T\boldsymbol{\Phi})\mathbf{w} - 2\mathbf{w}^T(\beta\boldsymbol{\Phi}^T)\mathbf{y} + \beta\mathbf{t}^T\mathbf{t} + \mathbf{w}^T\text{diag}(\boldsymbol{\alpha})\mathbf{w} = \\ & \mathbf{w}^T(\text{diag}(\boldsymbol{\alpha}) + \beta\boldsymbol{\Phi}^T\boldsymbol{\Phi})\mathbf{w} - 2\mathbf{w}^T(\beta\boldsymbol{\Phi}^T)\mathbf{t} = \\ & \mathbf{w}^T\boldsymbol{\Sigma}^{-1}\mathbf{w} - 2\mathbf{w}^T\boldsymbol{\Sigma}^{-1}(\beta\boldsymbol{\Sigma}\boldsymbol{\Phi}^T)\mathbf{t}. \end{aligned}$$

From which the result follows directly. End of Proof.

Important remark: When completing the square for a distribution, you can subtract, add or even ignore anything that does not depend on the random variable whose posterior you are calculating!

Or using Results for Linear Gaussian Models

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$

$$p(\mathbf{y} \mid \mathbf{x}) = \mathcal{N}(\mathbf{y} \mid \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$

□ For the above linear Gaussian model, the following hold:

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} \mid \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T)$$

$$p(\mathbf{x} \mid \mathbf{y}) = \mathcal{N}\left(\mathbf{x} \mid \left(\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A}\right)^{-1} \left(\boldsymbol{\Lambda} \boldsymbol{\mu} + \mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b})\right), \left(\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A}\right)^{-1}\right)$$

□ For our problem $\mathbf{x} \rightarrow \mathbf{w}$, $\boldsymbol{\mu} \rightarrow \mathbf{0}$, $\boldsymbol{\Lambda} \rightarrow \text{diag}(\boldsymbol{\alpha})$, $\mathbf{y} \rightarrow \mathbf{t}$, $\mathbf{A} \rightarrow \boldsymbol{\Phi}$, $\mathbf{b} \rightarrow \mathbf{0}$, $\mathbf{L} = \beta \mathbf{I}_N$

$$p(\mathbf{t} \mid \mathbf{w}) = \mathcal{N}_N(\mathbf{y} \mid \boldsymbol{\Phi} \mathbf{w}, \beta^{-1} \mathbf{I}_N), p(\mathbf{w}) = \mathcal{N}_M(\mathbf{w} \mid \text{diag}(\boldsymbol{\alpha})^{-1})$$

$$\text{and: } p(\mathbf{w} \mid \mathcal{D}, \beta, \boldsymbol{\alpha}) = \mathcal{N}\left(\mathbf{w} \mid \beta \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{t}, \underbrace{\left(\text{diag}(\boldsymbol{\alpha}) + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}\right)^{-1}}_{\boldsymbol{\Sigma}}\right)$$

RVM: Predictive Distribution

The **predictive** distribution is easily found as follows:

$$\begin{aligned} p(t|\mathbf{x}, \mathbf{X}, \mathbf{t}, \boldsymbol{\alpha}, \beta) &= \int p(t|\mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w}|\mathbf{X}, \mathbf{t}, \boldsymbol{\alpha}, \beta) d\mathbf{w} \\ &= \mathcal{N}(t|\mathbf{m}^T \boldsymbol{\phi}(\mathbf{x}), \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\Sigma} \boldsymbol{\phi}(\mathbf{x}) + \beta^{-1}). \end{aligned}$$

So, we have the *predictive mean* (\mathbf{m} being the posterior mean):

$$\mu(\mathbf{x}) = \mathbf{m}^T \boldsymbol{\phi}(\mathbf{x}),$$

and the predictive variance:

$$\sigma^2(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\Sigma} \boldsymbol{\phi}(\mathbf{x}) + \beta^{-1}.$$

Proof:

Complete the square....

End of proof.

RVM: Evidence Approximation

- We will select the hyper-parameters α and β by maximizing employing the *evidence approximation* (otherwise known as *maximizing the marginal likelihood*). We start by motivating it.
- Assume that we have a prior for α and β , $p(\alpha, \beta)$. Then the posterior of these parameters is:

$$p(\alpha, \beta | \mathcal{D}) \propto p(\mathcal{D} | \alpha, \beta) p(\alpha, \beta) = \int p(\mathcal{D} | \mathbf{w}, \beta) p(\mathbf{w} | \alpha) d\mathbf{w} p(\alpha, \beta).$$



Marginal Likelihood (or Evidence)

- Under the evidence approximation, we assume that:
 - $p(\alpha, \beta)$ is relatively flat.
 - $p(\mathcal{D} | \alpha, \beta)$ has very sharp maxima.

Then, we may obtain point estimates of the hyper-parameters by identifying these maxima.

RVM: Evidence

- To carry out the evidence approximation, we must obtain an analytic expression for the marginal likelihood. This is possible for RVM:

$$\mathcal{E}(\boldsymbol{\alpha}, \beta) = \log(p(\mathcal{D}|\boldsymbol{\alpha}, \beta)) = -\frac{1}{2} [N \log(2\pi) + \log|\mathbf{C}| + \mathbf{t}^T \mathbf{C}^{-1} \mathbf{t}],$$

where $\mathbf{C} \in \mathbf{R}^{N \times N}$, is given by:

$$\mathbf{C} = \beta^{-1} \mathbf{I}_N + \boldsymbol{\Phi} \text{diag}(\boldsymbol{\alpha})^{-1} \boldsymbol{\Phi}^T.$$

Proof:

This does require some more work that we avoided doing before. We start again from:

$$p(\mathcal{D}|\mathbf{w}, \beta) p(\mathbf{w}|\boldsymbol{\alpha}) = \mathcal{N}_N(\mathbf{t}|\boldsymbol{\Phi}\mathbf{w}, \beta^{-1} \mathbf{I}_N) \mathcal{N}_M(\mathbf{w}|\text{diag}(\boldsymbol{\alpha})^{-1}),$$

and we will complete the square making sure this time that we keep track of all terms that depend on \mathbf{t} .

RVM: Evidence

Proof (continuation):

Remember, that we want to integrate \mathbf{w} out of:

$$p(\mathcal{D}|\mathbf{w}, \beta)p(\mathbf{w}|\boldsymbol{\alpha}) = \mathcal{N}_n(\mathbf{t}|\boldsymbol{\Phi}\mathbf{w}, \beta^{-1}\mathbf{I}_N)\mathcal{N}_M(\mathbf{w}|\text{diag}(\boldsymbol{\alpha})^{-1}).$$

Again, we look at whatever is inside the exponential of the right-hand side and attempt to complete the square. This time, we already know the answer, but we need to keep track of \mathbf{t} explicitly:

$$\begin{aligned} & (\boldsymbol{\Phi}\mathbf{w} - \mathbf{t})^T(\beta\mathbf{I}_N)(\boldsymbol{\Phi}\mathbf{w} - \mathbf{t}) + \mathbf{w}^T\text{diag}(\boldsymbol{\alpha})\mathbf{w} = \\ & \mathbf{w}^T(\beta\boldsymbol{\Phi}^T\boldsymbol{\Phi})\mathbf{w} - 2\mathbf{w}^T(\beta\boldsymbol{\Phi}^T)\mathbf{t} + \beta\mathbf{t}^T\mathbf{t} + \mathbf{w}^T\text{diag}(\boldsymbol{\alpha})\mathbf{w} = \\ & \mathbf{w}^T(\text{diag}(\boldsymbol{\alpha}) + \beta\boldsymbol{\Phi}^T\boldsymbol{\Phi})\mathbf{w} - 2\mathbf{w}^T(\beta\boldsymbol{\Phi}^T)\mathbf{t} + \beta\mathbf{t}^T\mathbf{t} = \\ & \mathbf{w}^T\boldsymbol{\Sigma}^{-1}\mathbf{w} - 2\mathbf{w}^T\boldsymbol{\Sigma}^{-1}(\beta\boldsymbol{\Sigma}\boldsymbol{\Phi}^T)\mathbf{t} + \beta\mathbf{t}^T\mathbf{t} = \\ & \mathbf{w}^T\boldsymbol{\Sigma}^{-1}\mathbf{w} - 2\mathbf{w}^T\boldsymbol{\Sigma}^{-1}\mathbf{m} + \mathbf{m}^T\boldsymbol{\Sigma}^{-1}\mathbf{m} - \mathbf{m}^T\boldsymbol{\Sigma}^{-1}\mathbf{m} + \beta\mathbf{t}^T\mathbf{t} = \\ & (\mathbf{w} - \mathbf{m})^T\boldsymbol{\Sigma}^{-1}(\mathbf{w} - \mathbf{m}) + (\beta\mathbf{t}^T\mathbf{t} - \mathbf{m}^T\boldsymbol{\Sigma}^{-1}\mathbf{m}). \end{aligned}$$

RVM: Evidence

Proof (continuation):

We can now write:

$$p(\mathcal{D}|\mathbf{w}, \beta)p(\mathbf{w}|\boldsymbol{\alpha}) = \mathcal{N}_N(\mathbf{y}|\boldsymbol{\Phi}\mathbf{w}, \beta^{-1}\mathbf{I}_N)\mathcal{N}_M(\mathbf{w}|\text{diag}(\boldsymbol{\alpha})^{-1}) = \\ (2\pi)^{-\frac{N}{2}}\beta^{\frac{N}{2}}(2\pi)^{-\frac{M}{2}}|\text{diag}(\boldsymbol{\alpha})|^{\frac{1}{2}}\exp\left\{-\frac{1}{2}(\mathbf{w} - \mathbf{m})^T\boldsymbol{\Sigma}^{-1}(\mathbf{w} - \mathbf{m})\right\}$$

RVM: Two useful matrix identities

□ Woodbury Matrix Identity:

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}.$$

□ Matrix Determinant Lemma:

$$|A + UWV^T| = |W^{-1} + V^T A^{-1}U| |W| |A|.$$

You do not have to remember these, but you must be aware of their existence. Sometimes, they are the only route to analytic progress!

Or using Results for Linear Gaussian Models

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$

$$p(\mathbf{y} \mid \mathbf{x}) = \mathcal{N}(\mathbf{y} \mid \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$

□ For the above linear Gaussian model, the following hold:

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} \mid \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T)$$

$$p(\mathbf{x} \mid \mathbf{y}) = \mathcal{N}\left(\mathbf{x} \mid \left(\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A}\right)^{-1} \left(\boldsymbol{\Lambda} \boldsymbol{\mu} + \mathbf{A}^T \mathbf{L} (\mathbf{y} - \mathbf{b})\right), \left(\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A}\right)^{-1}\right)$$

□ For our problem $\mathbf{x} \rightarrow \mathbf{w}$, $\boldsymbol{\mu} \rightarrow \mathbf{0}$, $\boldsymbol{\Lambda} \rightarrow \text{diag}(\boldsymbol{\alpha})$, $\mathbf{y} \rightarrow \mathbf{t}$, $\mathbf{A} \rightarrow \boldsymbol{\Phi}$, $\mathbf{b} \rightarrow \mathbf{0}$, $\mathbf{L} = \beta \mathbf{I}_N$

$$p(\mathbf{t} \mid \mathbf{w}) = \mathcal{N}_N(\mathbf{y} \mid \boldsymbol{\Phi} \mathbf{w}, \beta^{-1} \mathbf{I}_N), p(\mathbf{w}) = \mathcal{N}_M(\mathbf{w} \mid \mathbf{0}, \text{diag}(\boldsymbol{\alpha})^{-1})$$

$$\text{and: } p(\mathbf{t} \mid \beta, \boldsymbol{\alpha}) = \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \beta^{-1} \mathbf{I}_N + \boldsymbol{\Phi} \text{diag}(\boldsymbol{\alpha})^{-1} \boldsymbol{\Phi}^T)$$

RVM: Evidence Approximation

- From [our earlier derivation](#):

$$p(\mathcal{D}|\boldsymbol{\alpha}, \beta) = (2\pi)^{-\frac{N}{2}} \beta^{\frac{N}{2}} |\text{diag}(\boldsymbol{\alpha})|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} \mathbf{t}^T (\beta \mathbf{I}_N - \beta^2 \boldsymbol{\Phi} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T) \mathbf{t} \right\} |\boldsymbol{\Sigma}|^{\frac{1}{2}}$$

$$\text{where } \boldsymbol{\Sigma} = (\text{diag}(\boldsymbol{\alpha}) + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1}, \quad \mathbf{m} = \beta \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{t}$$

- Noting that $\mathbf{t}^T \beta^2 \boldsymbol{\Phi} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{t} = \mathbf{m}^T \boldsymbol{\Sigma}^{-1} \mathbf{m}$, can simplify as:

$$p(\mathcal{D}|\boldsymbol{\alpha}, \beta) = (2\pi)^{-\frac{N}{2}} \beta^{\frac{N}{2}} |\text{diag}(\boldsymbol{\alpha})|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - \mathbf{m}^T \boldsymbol{\Sigma}^{-1} \mathbf{m}) \right\} |\boldsymbol{\Sigma}|^{\frac{1}{2}}$$

- Taking derivatives wrt α_i of $\ln p(\mathcal{D}|\boldsymbol{\alpha}, \beta)$ and using $\frac{\partial \ln |\boldsymbol{\Sigma}|}{\partial \alpha_i} =$

$$\text{Tr} \left(\boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\Sigma}}{\partial \alpha_i} \right) = -\text{Tr} \left(\boldsymbol{\Sigma} \frac{\partial \boldsymbol{\Sigma}^{-1}}{\partial \alpha_i} \right) = -\Sigma_{ii}, \quad \mathbf{m}^T \boldsymbol{\Sigma}^{-1} \frac{\partial \mathbf{m}}{\partial \alpha_i} =$$

$$-\mathbf{m}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma} \beta \frac{\partial \boldsymbol{\Sigma}^{-1}}{\partial \alpha_i} \boldsymbol{\Phi}^T \mathbf{t} = -m_i^2 \text{ and setting equal to 0:}$$

$$\frac{1}{2\alpha_i} + \frac{1}{2} m_i^2 - m_i^2 - \frac{1}{2} \Sigma_{ii} = 0 \rightarrow \alpha_i = \frac{1 - \alpha_i \Sigma_{ii}}{m_i^2} \equiv \frac{\gamma_i}{m_i^2}$$

RVM: Evidence Approximation

$$p(\mathcal{D}|\boldsymbol{\alpha}, \beta) = (2\pi)^{-\frac{N}{2}} \beta^{\frac{N}{2}} |\text{diag}(\boldsymbol{\alpha})|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - \mathbf{m}^T \boldsymbol{\Sigma}^{-1} \mathbf{m}) \right\} |\boldsymbol{\Sigma}|^{\frac{1}{2}}$$

□ Taking derivatives wrt β of $\ln p(\mathcal{D}|\boldsymbol{\alpha}, \beta)$ and using $\frac{\partial \ln |\boldsymbol{\Sigma}|}{\partial \beta} = \text{Tr} \left(\boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\Sigma}}{\partial \beta} \right) = -\text{Tr} \left(\boldsymbol{\Sigma} \frac{\partial \boldsymbol{\Sigma}^{-1}}{\partial \beta} \right) = -\text{Tr}(\boldsymbol{\Sigma} \boldsymbol{\Phi}^T \boldsymbol{\Phi})$, and $\mathbf{m} = \beta \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{t}$, setting equal to 0:

$$\frac{N}{2\beta} - \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{1}{2} \mathbf{m}^T \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{m} + \mathbf{m}^T \boldsymbol{\Phi}^T \mathbf{t} - \frac{1}{2} \text{tr}(\boldsymbol{\Sigma} \boldsymbol{\Phi}^T \boldsymbol{\Phi}) = 0 \rightarrow (\beta^{(new)})^{-1} = \frac{\|\mathbf{t} - \boldsymbol{\Phi} \mathbf{m}\|^2}{N - \sum \gamma_i}$$

□ Here, we used

$$\begin{aligned} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \boldsymbol{\Phi} &= \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \boldsymbol{\Phi} + \beta^{-1} \boldsymbol{\Sigma} \text{diag}(\boldsymbol{\alpha}) - \beta^{-1} \boldsymbol{\Sigma} \text{diag}(\boldsymbol{\alpha}) = \\ &= \boldsymbol{\Sigma} (\text{diag}(\boldsymbol{\alpha}) + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}) \beta^{-1} - \beta^{-1} \boldsymbol{\Sigma} \text{diag}(\boldsymbol{\alpha}) \\ &= (\text{diag}(\boldsymbol{\alpha}) + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} (\text{diag}(\boldsymbol{\alpha}) + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}) \beta^{-1} - \beta^{-1} \boldsymbol{\Sigma} \text{diag}(\boldsymbol{\alpha}) \\ &= \beta^{-1} (\mathbf{I} - \text{diag}(\boldsymbol{\alpha}) \boldsymbol{\Sigma}) \end{aligned}$$

and thus $\text{tr}(\boldsymbol{\Sigma} \boldsymbol{\Phi}^T \boldsymbol{\Phi}) = \beta^{-1} \sum \gamma_i$

RVM: Evidence Approximation

$$\alpha_i = \frac{1 - \alpha_i \Sigma_{ii}}{m_i^2} \equiv \frac{\gamma_i}{m_i^2}, \quad (\beta^{(new)})^{-1} = \frac{\|t - \Phi m\|^2}{N - \sum \gamma_i}$$

□ Learning therefore proceeds by

- ✓ choosing initial values for α and β ,
- ✓ evaluating the mean and covariance of the posterior using $\Sigma = (\text{diag}(\alpha) + \beta \Phi^T \Phi)^{-1}$, $m = \beta \Sigma \Phi^T t$, respectively, and then
- ✓ alternately re-estimating the hyperparameters, using the Eqs. above and
- ✓ iterating until a suitable convergence criterion is satisfied.

□ The second alternative approach is to use the EM algorithm. These two approaches to finding the values of the hyperparameters that maximize the evidence are formally equivalent. Numerically, however, it is found that the direct optimization approach gives somewhat faster convergence (Tipping, 2001).

RVM: Fast training algorithm

- ❑ We will develop a **fast algorithm** for **maximizing the evidence**.
- ❑ We follow closely the developments in [*Tipping, M. E. and A. C. Faul \(2003\). Fast marginal likelihood maximisation for sparse Bayesian models. In C. M. Bishop and B. J. Frey \(Eds.\), Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics, Key West, FL, Jan 3-6.*](#)
- ❑ There are some typos [in the paper](#).
- ❑ We will achieve this by asking the following question:

Keeping β fixed and given some values for α , what is the action on a single basis function that gives the maximum increase in evidence?

This will require a very thorough examination of the evidence function.

- ❑ A Matlab implementation for the one dimensional case can be found [here](#).

RVM: Splitting the evidence

□ Some notation:

- Subscript “ $-i$ ”: Wherever we see this applied to a vector or matrix, it gives the corresponding vector or matrix with the influence of the i –th basis function removed.
- $\alpha_{-i} = (\alpha_1, \dots, \alpha_{i-1}, \alpha_{i+1}, \dots, \alpha_m)$: All the hyper-parameters except the ones that pertain to basis function i .
- Φ_{-i} : The design matrix with the column corresponding to basis function i removed.
- $C_{-i} = \beta^{-1}I_N + \Phi_{-i}\text{diag}(\alpha_{-i})^{-1}\Phi_{-i}^T$: The matrix C , without the influence of basis function i .
- $\mathcal{E}(\alpha_{-i}, \beta) = -\frac{1}{2} [n \log(2\pi) + \log|C_{-i}| + \mathbf{y}^T C_{-i}^{-1} \mathbf{y}]$: The evidence when the i –th basis function has been removed.

RVM: Splitting the evidence

Given any basis function $i \in \{1, \dots, M\}$, the evidence can be split as:

$$\mathcal{E}(\boldsymbol{\alpha}, \beta) = \mathcal{E}(\boldsymbol{\alpha}_{-i}, \beta) + \epsilon(\alpha_i),$$

where

$$\epsilon(\alpha_i) = \frac{1}{2} \left(\log \alpha_i - \log(\alpha_i + s_i) + \frac{h_i^2}{\alpha_i + s_i} \right),$$

with

$$s_i = \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i \quad \text{and} \quad h_i = \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \mathbf{t}.$$

Proof:

Remark:

Everything that relates to basis function i is here!

We start by splitting \mathbf{C} :

$$\begin{aligned} \mathbf{C} &= \beta^{-1} \mathbf{I}_N + \boldsymbol{\Phi} \text{diag}(\boldsymbol{\alpha})^{-1} \boldsymbol{\Phi}^T \\ &= \beta^{-1} \mathbf{I}_N + (\boldsymbol{\Phi}_{-i} \quad \boldsymbol{\varphi}_i) \begin{pmatrix} \text{diag}(\boldsymbol{\alpha}_{-i})^{-1} & \mathbf{0}_{(M-1) \times 1} \\ \mathbf{0}_{1 \times (M-1)} & \alpha_i^{-1} \end{pmatrix} \begin{pmatrix} \boldsymbol{\Phi}_{-i}^T \\ \boldsymbol{\varphi}_i^T \end{pmatrix} \\ &= \beta^{-1} \mathbf{I}_N + \boldsymbol{\Phi}_{-i} \text{diag}(\boldsymbol{\alpha}_{-i})^{-1} \boldsymbol{\Phi}_{-i}^T + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^T \\ &= \mathbf{C}_{-i} + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^T. \end{aligned}$$

RVM: Splitting the evidence

Proof (continuation):

Remember that the evidence is:

$$\mathcal{E}(\boldsymbol{\alpha}, \beta) = -\frac{1}{2} [N \log(2\pi) + \log|\mathbf{C}| + \mathbf{t}^T \mathbf{C}^{-1} \mathbf{t}].$$

We will only make use of this:

$$\mathbf{C} = \mathbf{C}_{-i} + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^T,$$

and some identities.

From the Woodbury identity, we get:

$$\begin{aligned} & (\mathbf{A} + \mathbf{UCV})^{-1} \\ &= \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{U} (\mathbf{C}^{-1} + \mathbf{VA}^{-1} \mathbf{U})^{-1} \mathbf{VA}^{-1}. \end{aligned}$$

$$\begin{aligned} \mathbf{C}^{-1} &= \mathbf{C}_{-i}^{-1} - \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i (\alpha_i + \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i)^{-1} \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \\ &= \mathbf{C}_{-i}^{-1} - \frac{\mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1}}{\alpha_i + s_i}. \end{aligned}$$

And we get:

$$\mathbf{t}^T \mathbf{C}^{-1} \mathbf{t} = \mathbf{t}^T \mathbf{C}_{-i}^{-1} \mathbf{t} - \frac{(\mathbf{t}^T \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i)(\boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \mathbf{t})}{\alpha_i + s_i} = \mathbf{t}^T \mathbf{C}_{-i}^{-1} \mathbf{t} + \frac{h_i^2}{\alpha_i + s_i}.$$

RVM: Splitting the evidence

Proof (continuation):

Finally, we look at $|\mathbf{C}|$. Remember:

$$\mathbf{C} = \mathbf{C}_{-i} + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^T,$$

Using the matrix determinant lemma:

$$|\mathbf{A} + \mathbf{U}\mathbf{W}\mathbf{V}^T| = |\mathbf{W}^{-1} + \mathbf{V}^T \mathbf{A}^{-1} \mathbf{U}| |\mathbf{W}| |\mathbf{A}|$$

$$\begin{aligned} |\mathbf{C}| &= |\alpha_i + \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i| |\alpha_i^{-1}| |\mathbf{C}_{-i}| \\ &= \frac{|\mathbf{C}_{-i}| (\alpha_i + s_i)}{\alpha_i}. \end{aligned}$$

Or

$$\log |\mathbf{C}| = \log |\mathbf{C}_{-i}| - \log \alpha_i + \log(\alpha_i + s_i).$$

RVM: Splitting the evidence

Proof (continuation):

Now we combine everything. Remember that:

$$\mathbf{t}^T \mathbf{C}^{-1} \mathbf{t} = \mathbf{t}^T \mathbf{C}_{-i}^{-1} \mathbf{t} + \frac{h_i^2}{\alpha_i + s_i},$$

and

$$\log |\mathbf{C}| = \log |\mathbf{C}_{-i}| - \log \alpha_i + \log(\alpha_i + s_i).$$

Therefore, we have for the evidence:

$$\begin{aligned} \mathcal{E}(\boldsymbol{\alpha}, \beta) &= -\frac{1}{2} [N \log(2\pi) + \log |\mathbf{C}| + \mathbf{t}^T \mathbf{C}^{-1} \mathbf{t}] \\ &= -\frac{1}{2} \left(N \log(2\pi) + \log |\mathbf{C}_{-i}| + \mathbf{t}^T \mathbf{C}_{-i}^{-1} \mathbf{t} + \frac{h_i^2}{\alpha_i + s_i} - \log \alpha_i + \log(\alpha_i + s_i) \right) \\ &= \mathcal{E}(\boldsymbol{\alpha}_{-i}, \beta) + \epsilon(\alpha_i). \end{aligned}$$

End of proof.

RVM: Analysis of the evidence

We have shown that for any basis function $i \in \{1, \dots, M\}$, the evidence can be split as:

$$\mathcal{E}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \mathcal{E}(\boldsymbol{\alpha}_{-i}, \boldsymbol{\beta}) + \epsilon(\alpha_i),$$

where

$$\epsilon(\alpha_i) = \frac{1}{2} \left(\log \alpha_i - \log(\alpha_i + s_i) + \frac{h_i^2}{\alpha_i + s_i} \right).$$

We will now analyze the $\epsilon(\alpha_i)$ term in more detail.

We compute its first derivative:

$$\frac{\partial \epsilon(\alpha_i)}{\partial \alpha_i} = \frac{s_i^2 - \alpha_i \theta_i}{2\alpha_i(\alpha_i + s_i)^2},$$

where

$$\theta_i = h_i^2 - s_i.$$

RVM: Analysis of the evidence

Setting

$$\frac{\partial \epsilon(\alpha_i)}{\partial \alpha_i} = \frac{s_i^2 - \alpha_i \theta_i}{2\alpha_i(\alpha_i + s_i)^2}$$

equal to zero, we have two possibilities:

□ If $\theta_i > 0$, then

$$\alpha_i^* = \frac{s_i^2}{\theta_i},$$

is the unique maximum of $\epsilon(\alpha_i)$ (calculate the 2nd derivative to be convinced that it is indeed a maximum).

□ If $\theta_i \leq 0$, then

$$\alpha_i^* = +\infty, \quad \text{It is monotonic}$$

Is the unique maximum of $\epsilon(\alpha_i)$ (convince yourselves that its derivative remains positive for all α_i). This means that removing the corresponding basis function from the model, increases the evidence.

RVM: Possible actions

The previous observations translate into three possible actions for each basis function i :

□ **ADD**, if $\theta_i > 0$ and $\alpha_i = +\infty$:

$$\alpha_i^{\text{new}} = \frac{s_i^2}{\theta_i},$$
$$\Delta\mathcal{E} = \epsilon(\alpha_i^{\text{new}}).$$

□ **RE-ESTIMATE**, if $\theta_i > 0$ and $\alpha_i < +\infty$:

$$\alpha_i^{\text{new}} = \frac{s_i^2}{\theta_i},$$
$$\Delta\mathcal{E} = \epsilon(\alpha_i^{\text{new}}) - \epsilon(\alpha_i).$$

□ **DELETE**, if $\theta_i \leq 0$ and $\alpha_i < +\infty$:

$$\alpha_i^{\text{new}} = +\infty,$$
$$\Delta\mathcal{E} = -\epsilon(\alpha_i).$$

RVM: Updating the noise

Now, we assume that all the α 's are kept fixed and we attempt to maximize the evidence with respect to the inverse noise β .

Taking the derivative of the evidence with respect to β , you get:

$$\beta^{\text{new}} = \frac{N - M - \sum_i \alpha_i \Sigma_{ii}}{\| \mathbf{t} - \mathbf{\Phi} \boldsymbol{\mu} \|^2}.$$

This is a lengthy calculation. You are going to need the following formulas for a general matrix \mathbf{A} :

$$\frac{d\mathbf{A}^{-1}}{dt} = -\mathbf{A}^{-1} \frac{d\mathbf{A}}{dt} \mathbf{A}^{-1},$$

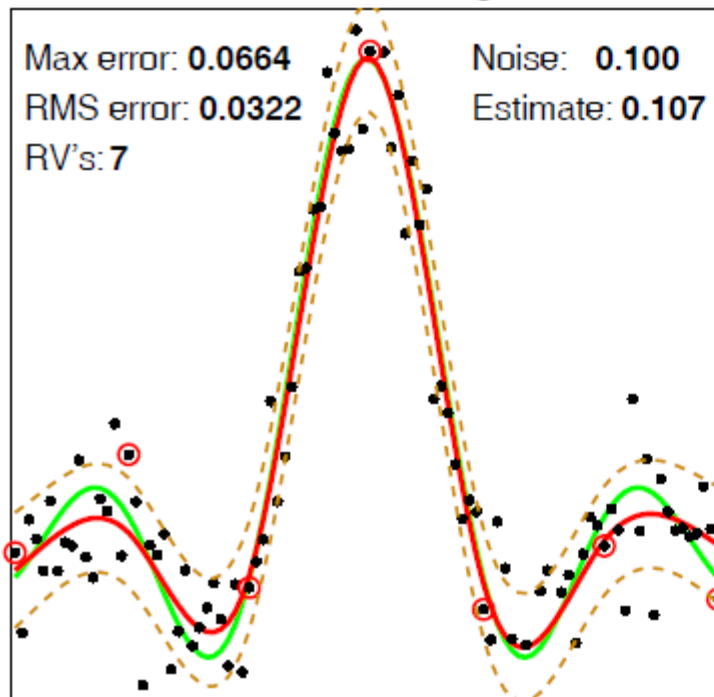
$$\frac{d|\mathbf{A}|}{dt} = |\mathbf{A}| \text{tr} \left[\mathbf{A}^{-1} \frac{d\mathbf{A}}{dt} \right].$$

RVM: Complete Algorithm

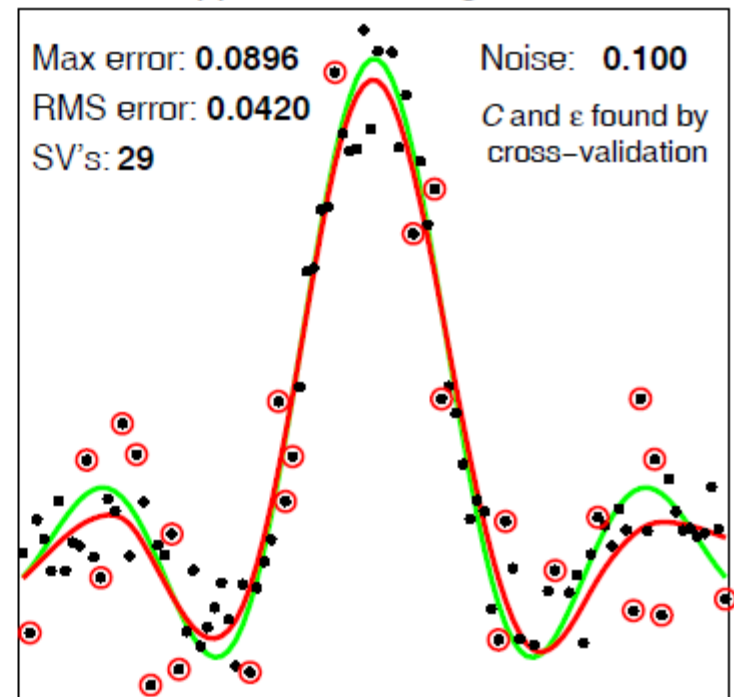
1. Initialize β and α (Some of them might be $+\infty$. The corresponding basis functions are not in the model).
2. Compute: Σ , \mathbf{m} only for the basis functions in the model but **all** the s_i and h_i 's.
3. For $i = 1, \dots, M$, compute $\theta_i = h_i^2 - s_i$ and find the basis function i^* and the corresponding action \mathcal{A}_{i^*} that gives the maximum increase in evidence.
4. Perform the best action.
5. Re-estimate the noise β .
6. Re-compute (or update) Σ , \mathbf{m} , s_i and h_i .
7. If converged, stop. Otherwise go to 3.

A simple example

Relevance Vector Regression



Support Vector Regression



Fitting noisy samples from $f(x) = \frac{\sin(x)}{x}$. The basis functions are Gaussian kernels centered on the data points:

$$\phi_i(x) = k(x, x^{(i)}) = \exp\left(-\lambda(x - x^{(i)})^2\right).$$

Taken [from this](#) very nice [tutorial](#). The comparison is with [SVM](#) (a non-Bayesian competitor).

RVM: Fast updates

An efficient implementation of the RVM algorithm, requires a fast way to update the various statistics.

First, the mean and the co-variance matrix should only be build out of the **relevant** basis functions:

$$\begin{aligned}\Sigma_r &= (\text{diag}(\alpha_r) + \beta \Phi_r^T \Phi_r)^{-1}, \\ \mathbf{m}_r &= \beta \Sigma_r \Phi_r^T \mathbf{y},\end{aligned}$$

where the subscript “r” denotes that the corresponding quantity corresponds only to the basis functions in the model (the relevant basis functions), i.e. the ones with $\alpha_i < +\infty$. The other basis functions are simply irrelevant...

This requires some book-keeping. That is, at any iteration you must keep track of the set of relevant basis functions:

$$\mathcal{R} = \{j \in \{1, \dots, m\}: \alpha_j < +\infty\}.$$

RVM: Fast updates

Secondly, we come to the statistics:

$$s_i = \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i \quad \text{and} \quad h_i = \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \mathbf{t}.$$

These are required for **all** basis functions, so that the change in evidence can be computed. It is more convenient to keep track of the following statistics instead:

$$S_i = \boldsymbol{\varphi}_i^T \mathbf{C}^{-1} \boldsymbol{\varphi}_i \quad \text{and} \quad H_i = \boldsymbol{\varphi}_i^T \mathbf{C}^{-1} \mathbf{t}.$$

It is trivial to show that they connect with the original by:

$$s_i = \frac{\alpha_i S_i}{\alpha_i - S_i} \quad \text{and} \quad h_i = \frac{\alpha_i H_i}{\alpha_i - S_i}.$$

Using the Woodbury identity once more, you can find an easy way to compute the 'capital statistics':

$$\begin{aligned} S_i &= \beta \|\boldsymbol{\varphi}_i\|^2 - \beta^2 \boldsymbol{\varphi}_i^T \boldsymbol{\Phi}_r \boldsymbol{\Sigma}_r \boldsymbol{\Phi}_r^T \boldsymbol{\varphi}_i, \\ H_i &= \beta \boldsymbol{\varphi}_i^T \mathbf{t} - \beta^2 \boldsymbol{\varphi}_i^T \boldsymbol{\Phi}_r \boldsymbol{\Sigma}_r \boldsymbol{\Phi}_r^T \mathbf{t}. \end{aligned}$$

RVM: Fast updates

Finally, each action can be carried out without computing the statistics from scratch! You can derive them by repeatedly applying the Woodbury identity. It will take a couple of pages...

Notation:

- M_r : current number of relevant basis functions
- $i \in \{1, \dots, M\}$: index of basis function that is currently being updated.
- $j \in \{1, \dots, M_r\}$: the index within the **relevant** basis set that corresponds to i .
- $t \in \{1, \dots, M\}$: index that runs over all basis functions (in or out of the model).

RVM: Fast updates

ADD basis function i :

$$\begin{aligned}\Delta\mathcal{E}^{\text{new}} &= \frac{1}{2} \left(\frac{Q_i^2 - S_i}{S_i} + \log \frac{S_i}{Q_i^2} \right), \\ \Sigma_r^{\text{new}} &= \begin{pmatrix} \Sigma_r + \beta^2 \Sigma_{ii} \Sigma_r \Phi_r^T \varphi_i \varphi_i^T \Phi_r \Sigma_r & -\beta \Sigma_{ii} \Sigma_r \Phi_r^T \varphi_i \\ -\beta \Sigma_{ii} \varphi_i^T \Phi_r \Sigma_r & \Sigma_{ii} \end{pmatrix}, \\ \mathbf{m}_r^{\text{new}} &= \begin{pmatrix} \mathbf{m}_r - \mu_i \beta \Sigma_{ii} \Sigma_r \Phi_r^T \varphi_i \\ m_i \end{pmatrix}, \\ S_t^{\text{new}} &= S_t - \Sigma_{ii} (\beta \varphi_t^T \mathbf{e}_i)^2, \\ H_t^{\text{new}} &= H_t - \mu_i \beta \varphi_t^T \mathbf{e}_i,\end{aligned}$$

where

$$\begin{aligned}\Sigma_{ii} &= (\alpha_i^{\text{new}} - S_i)^{-1}, \\ \mu_i &= \Sigma_{ii} H_i,\end{aligned}$$

and

$$\mathbf{e}_i = \varphi_i - \beta \Phi_r \Sigma_r \Phi_r^T \varphi_i.$$

Finally, m_r is increased by one.

CAUTION: There is a typo in Tipping and Faul.

RVM: Fast updates

DELETE basis function j :

$$\begin{aligned}\Delta \mathcal{E}^{\text{new}} &= \frac{1}{2} \left(\frac{Q_i^2 - S_i}{S_i} - \log \left(1 - \frac{S_i}{\alpha_i} \right) \right), \\ \boldsymbol{\Sigma}_r^{\text{new}} &= \boldsymbol{\Sigma}_r - \Sigma_{jj}^{-1} \boldsymbol{\Sigma}_j \boldsymbol{\Sigma}_j^T, \\ \mathbf{m}_r^{\text{new}} &= \mathbf{m}_r - \mu_j \Sigma_{jj}^{-1} \boldsymbol{\Sigma}_j, \\ S_t^{\text{new}} &= S_t + \Sigma_{jj}^{-1} (\beta \boldsymbol{\Sigma}_j^T \boldsymbol{\Phi}_r^T \boldsymbol{\varphi}_t)^2, \\ H_t^{\text{new}} &= H_t + m_j \Sigma_{jj}^{-1} \beta \boldsymbol{\Sigma}_j^T \boldsymbol{\Phi}_r^T \boldsymbol{\varphi}_t,\end{aligned}$$

where Σ_{jj} is the (j, j) element and $\boldsymbol{\Sigma}_j$ the j -th column of $\boldsymbol{\Sigma}_r$, resp., and m_j is the j -th element of \mathbf{m}_r .

After completing this step the corresponding rows and columns of the mean and the covariance must be removed.

The indices of the relevant vectors must be updated accordingly.

RVM: Fast updates

RE-ESTIMATE basis function j :

$$\Delta \mathcal{E}^{\text{new}} = \frac{1}{2} \left(\frac{Q_i^2}{S_i + \left((\alpha_i^{\text{new}})^{-1} - \alpha_i^{-1} \right)^{-1}} - \log \left(1 - S_i \left((\alpha_i^{\text{new}})^{-1} - \alpha_i^{-1} \right) \right) \right),$$

$$\mathbf{\Sigma}_r^{\text{new}} = \mathbf{\Sigma}_r - \kappa_j \mathbf{\Sigma}_j \mathbf{\Sigma}_j^T,$$

$$\mathbf{m}_r^{\text{new}} = \mathbf{m}_r - \kappa_j m_j \mathbf{\Sigma}_j,$$

$$S_t^{\text{new}} = S_t + \kappa_j \left(\beta \mathbf{\Sigma}_j^T \mathbf{\Phi}_r^T \boldsymbol{\varphi}_t \right)^2,$$

$$H_t^{\text{new}} = H_t + \kappa_j m_j \beta \mathbf{\Sigma}_j^T \mathbf{\Phi}_r^T \boldsymbol{\varphi}_t,$$

$\mathbf{\Sigma}_j$ the j -th column of $\mathbf{\Sigma}_r$, resp., and μ_j is the j -th element of \mathbf{m}_r and

$$\kappa_j = \left(\Sigma_{jj} + \left((\alpha_i^{\text{new}})^{-1} - \alpha_i^{-1} \right)^{-1} \right)^{-1},$$

where Σ_{jj} is the (j, j) element.

RVM: Extension to multiple outputs

- We assume that we observe the data:

$$\mathcal{D} = \{\mathbf{x}^{(i)}, \mathbf{t}^{(i)} = \mathbf{f}(\mathbf{x}^{(i)})\}_{i=1}^n,$$

and that we wish to learn the function $\mathbf{f}(\cdot) \in \mathbf{R}^q$.

- Pick a set of M basis functions $\boldsymbol{\phi}: \mathbf{R}^k \rightarrow \mathbf{R}^M$:

$$\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})).$$

- We will approximate $\mathbf{f}(\cdot)$ by:

$$\hat{\mathbf{f}}(\mathbf{x}; \mathbf{W}) = \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}) + \boldsymbol{\epsilon},$$

where $\boldsymbol{\epsilon}$ is a random vector representing noise and $\mathbf{W} \in \mathbf{R}^{M \times q}$.

- For simplicity, we choose to work with Gaussian noise with zero mean and inverse variance β , so the **likelihood** is:

$$p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) = \mathcal{N}_q(\mathbf{t}|\mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}), \beta^{-1}).$$

RVM: Extension to multiple outputs

- The prior on the weights:

$$p(\mathbf{W}|\boldsymbol{\alpha}) = \prod_{j=1}^M p(\mathbf{w}_j|\alpha_j),$$

with $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_M)$ and:

$$p(\mathbf{w}_j|\alpha_j) = \mathcal{N}_q(\mathbf{w}_j|0, \alpha_j^{-1}).$$

- The posterior of the weights is:

$$p(\mathbf{W}|\mathcal{D}, \beta, \boldsymbol{\alpha}) = \mathcal{N}(\mathbf{W}|\mathbf{M}, \boldsymbol{\Sigma}),$$

where:

$$\begin{aligned}\boldsymbol{\Sigma} &= (\text{diag}(\boldsymbol{\alpha}) + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1}, \\ \mathbf{M} &= \beta \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{Y},\end{aligned}$$

with $\mathbf{Y} = (\mathbf{t}^{(1)}, \dots, \mathbf{t}^{(N)}) \in \mathbf{R}^{N \times q}$ is the vector of observations, $\boldsymbol{\Phi} \in \mathbf{R}^{N \times M}$ is the design matrix:

RVM: Predictive distribution

□ The **predictive** distribution is easily found as follows:

$$\begin{aligned} p(\mathbf{t}|\mathbf{x}, \alpha, \beta) &= \int p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) p(\mathbf{W}|\alpha) d\mathbf{W} \\ &= \mathcal{N}(\mathbf{y} | \mathbf{M}^T \boldsymbol{\phi}(\mathbf{x}), \boldsymbol{\varphi}(\mathbf{x})^T \boldsymbol{\Sigma} \boldsymbol{\varphi}(\mathbf{x}) + \beta^{-1}). \end{aligned}$$

So, we have the *predictive mean*:

$$\boldsymbol{\mu}(\mathbf{x}) = \mathbf{M}^T \boldsymbol{\phi}(\mathbf{x}),$$

and the predictive variance:

$$\sigma^2(\mathbf{x}) = \boldsymbol{\varphi}(\mathbf{x})^T \boldsymbol{\Sigma} \boldsymbol{\varphi}(\mathbf{x}) + \beta^{-1}.$$

Proof:

Complete the square....

End of proof.

RVM: Extension to multiple outputs

- **REMARK:** You might have to work with a scaled version of the data so that the prior assumptions are closer to reality. You may attempt to fit the function:

$$\mathbf{g}(\mathbf{x}) = \mathbf{L}_{\text{obs}}^{-1}(\mathbf{f}(\mathbf{x}) - \mathbf{m}_{\text{obs}}),$$

where \mathbf{m}_{obs} is the *empirical mean of the data*:

$$\mathbf{m}_{\text{obs}} = \frac{1}{n} \sum_{i=1}^N \mathbf{t}^{(i)},$$

and $\mathbf{L}_{\text{obs}} \in \mathbf{R}^{q \times q}$ is the Cholesky decomposition of the *empirical covariance matrix of the data*:

$$\mathbf{C}_{\text{obs}} = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{t}^{(i)} - \mathbf{m}_{\text{obs}})(\mathbf{t}^{(i)} - \mathbf{m}_{\text{obs}})^T.$$

If the empirical covariance is rank-deficient, you may use either an incomplete Cholesky or the eigen-decomposition. This is essentially a sort of output dimensionality reduction.

RVM: Extension to multiple outputs

- We proceed again, by maximizing the total evidence.
- Working out the details, it is easy to show that the total evidence is the sum of the evidence of each output:

$$\mathcal{E}(\boldsymbol{\alpha}, \beta | \mathcal{D}) = \sum_{j=1}^q \mathcal{E}(\boldsymbol{\alpha}, \beta | \mathcal{D}_j),$$

where $\mathcal{D}_j = \left\{ \left(\mathbf{x}^{(i)}, t_j^{(i)} \right) \right\}_{i=1}^N$.

- The algorithmic details remain essentially the same. The only difference is that:

$$H_i = \frac{1}{q} \sum_{j=1}^q H_{ji},$$

$$H_{ji} = \beta \boldsymbol{\varphi}_i^T \mathbf{t}_j - \beta^2 \boldsymbol{\varphi}_i^T \boldsymbol{\Phi}_r \boldsymbol{\Sigma}_r \boldsymbol{\Phi}_r^T \mathbf{t}_j.$$

RVM: Numerical stability concerns

- ❑ There are several numerical instabilities that should concern you.
- ❑ This is very common.
- ❑ For a numerically stable way to implement RVM using the Generalized Singular Value Decomposition look at Appendix D of [I. Bilonis and N. Zabarar. *Multidimensional adaptive relevance vector machines for uncertainty quantification*. SIAM Journal for Scientific Computing, Vol. 34, No. 6, pp. B881–B908 2012.](#)

RVM For Classification

- We can extend the relevance vector machine framework to classification problems by applying the ARD prior over weights to a probabilistic linear classification model.
- Consider two-class problems with a binary target variable $t \in \{0, 1\}$.
- The model now takes the form of a linear combination of basis functions transformed by a logistic sigmoid function

$$y(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}))$$

- If we introduce a Gaussian prior over \mathbf{w} , then we obtain the model that has been considered in the logistic regression lecture.

RVM For Classification

- In the RVM, this model uses the ARD prior

$$p(w_j | \alpha_j) = \mathcal{N}(w_j | 0, \alpha_j^{-1}).$$

- We can no longer integrate analytically over \mathbf{w} . Here we follow Tipping (2001) and use the Laplace approximation.
- We begin by initializing α . For this given value of α , we then build a Gaussian approximation to the posterior distribution and thereby obtain an approximation to the marginal likelihood ($\mathbf{A} = \text{diag}(\alpha_i)$)

$$\begin{aligned} \ln p(\mathbf{w} | \mathbf{t}, \alpha) &= \ln \{p(\mathbf{t} | \mathbf{w}) p(\mathbf{w} | \alpha)\} - \ln p(\mathbf{t} | \alpha) \\ &= \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln (1 - y_n)\} - \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \text{const.} \end{aligned}$$

- Maximization of this approximate marginal likelihood then leads to a re-estimated value for α , and the process is repeated until convergence.

RVM For Classification

- This can be done using iterative reweighted least squares (IRLS).
- We need the gradient and Hessian of the log posterior distribution

$$\begin{aligned}\nabla \ln p(\mathbf{w}|\mathbf{t}, \alpha) &= \mathbf{\Phi}^T(\mathbf{t} - \mathbf{y}) - \mathbf{A}\mathbf{w} \\ \nabla \nabla \ln p(\mathbf{w}|\mathbf{t}, \alpha) &= -(\mathbf{\Phi}^T \mathbf{B} \mathbf{\Phi} + \mathbf{A})\end{aligned}$$

- \mathbf{B} is an $N \times N$ diagonal matrix with elements $b_n = y_n(1 - y_n)$, the vector $\mathbf{y} = (y_1, \dots, y_N)^T$, and $\mathbf{\Phi}$ is the design matrix with elements $\Phi_{ni} = \phi_i(\mathbf{x}_n)$.
- At convergence of the IRLS algorithm, the negative Hessian represents the inverse covariance matrix for the Gaussian approximation to the posterior distribution.
- The mean & covariance of the Laplace approximation take the form

$$\begin{aligned}\mathbf{w} &= \mathbf{A}^{-1} \mathbf{\Phi}^T(\mathbf{t} - \mathbf{y}) \\ \mathbf{\Sigma} &= (\mathbf{\Phi}^T \mathbf{B} \mathbf{\Phi} + \mathbf{A})^{-1}\end{aligned}$$

RVM For Classification

- The marginal likelihood can be approximated as:

$$p(\mathbf{t}|\boldsymbol{\alpha}) = \int p(\mathbf{t}|\mathbf{w})p(\mathbf{w}|\boldsymbol{\alpha})d\mathbf{w} \approx p(\mathbf{t}|\mathbf{w}^*)p(\mathbf{w}^*|\boldsymbol{\alpha})(2\pi)^{M/2}|\boldsymbol{\Sigma}|^{1/2}$$

- If we substitute for $p(\mathbf{t}|\mathbf{w}^*)$ and $p(\mathbf{w}^*|\boldsymbol{\alpha})$ and then set the derivative wrt α_i equal to zero:

$$-\frac{1}{2}(w_i^*)^2 + \frac{1}{2\alpha_i} - \frac{1}{2}\Sigma_{ii} = 0$$

- Defining $\gamma_i = 1 - \alpha_i\Sigma_{ii}$, gives:

$$\alpha_i^{new} = \frac{\gamma_i}{(w_i^*)^2}$$

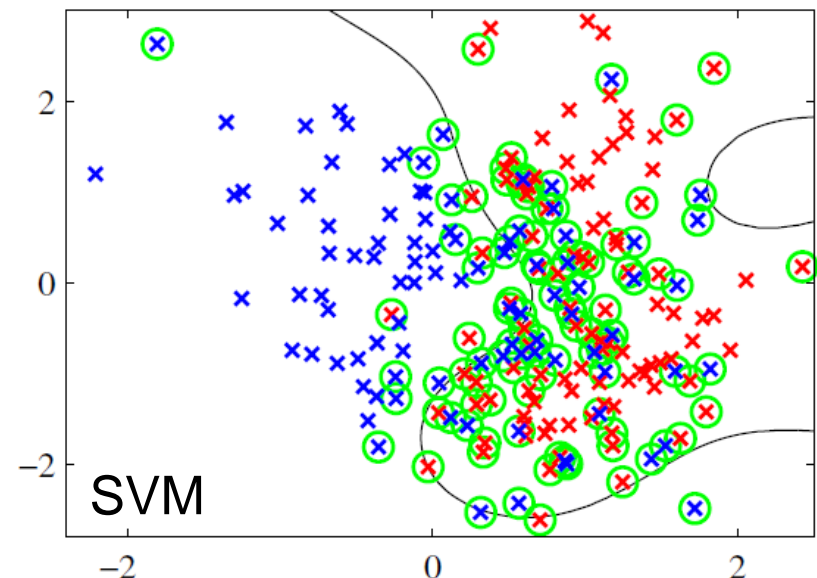
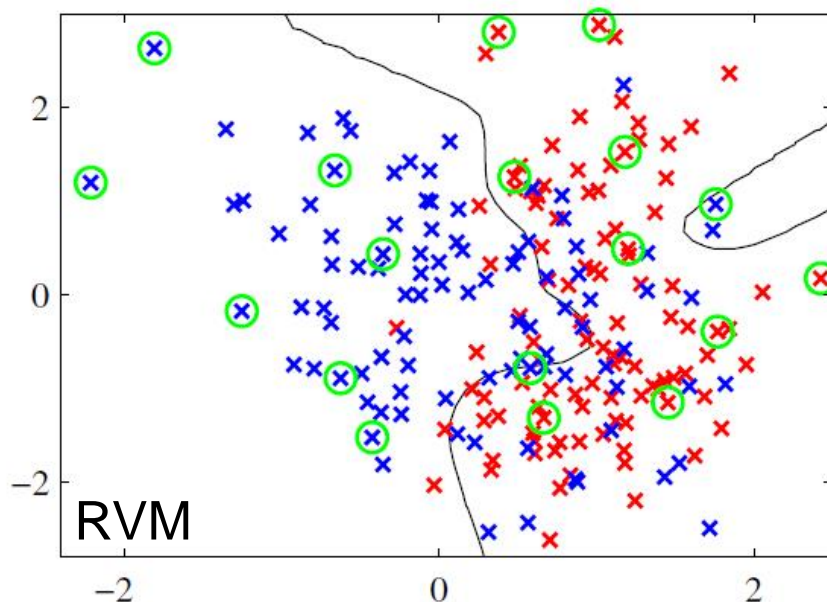
- With $\hat{\mathbf{t}} = \boldsymbol{\Phi}\mathbf{w}^* + \mathbf{B}^{-1}(\mathbf{t} - \mathbf{y})$, we can write

$$\ln p(\mathbf{t}|\boldsymbol{\alpha}) = -\frac{1}{2}\{N\ln(2\pi) + \ln|\mathbf{C}| + (\hat{\mathbf{t}})^T \mathbf{C}^{-1}\hat{\mathbf{t}}\}, \mathbf{C} = \mathbf{B} + \boldsymbol{\Phi}\mathbf{A}\boldsymbol{\Phi}^T$$

- This is the same form as for regression. Can apply the same sparsity analysis and obtain the same fast learning algorithm in which we optimize a single α_i at each step.

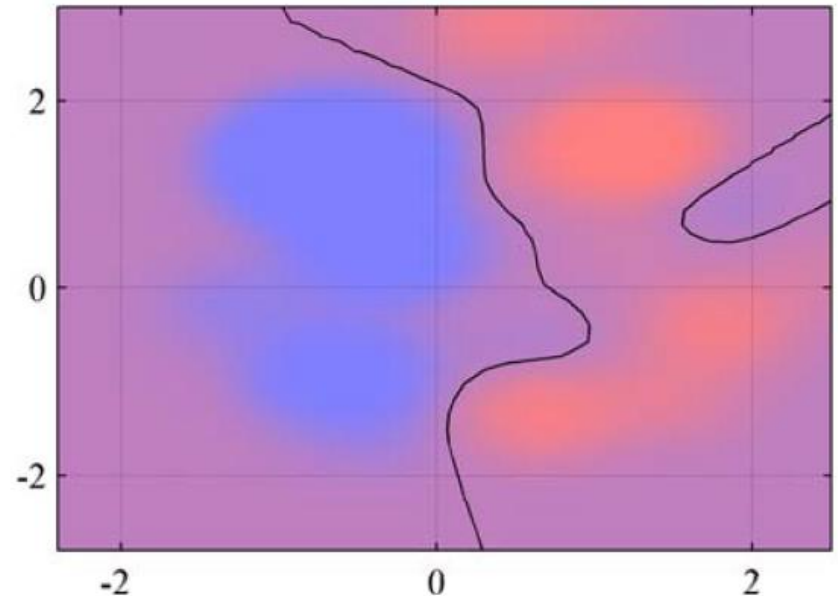
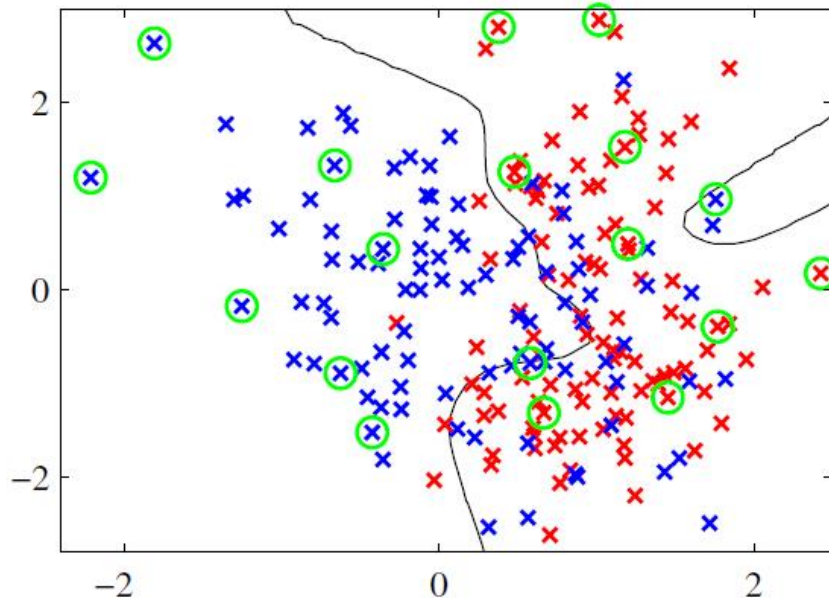
RVM For Classification

- Example of the RVM applied to a synthetic data set, in which the left-hand plot shows the decision boundary and the data points, with the relevance vectors indicated by circles.
- Comparing with the SVM results it is seen that the RVM gives a much sparser model.



RVM For Classification

- The right-hand plot shows the posterior probability given by the RVM output in which the proportion of red (blue) ink indicates the probability of that point belonging to the red (blue) class.



RVM For Classification

- So far, we have considered the RVM for binary classification problems. For $K > 2$ classes,

$$y_k(\mathbf{x}) = \frac{\exp(\alpha_k)}{\sum_j \exp(\alpha_k)}$$

- The log likelihood function is then given by

$$\ln p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$$

where the target values t_{nk} have a 1-of- K coding for each data point n , and \mathbf{T} is a matrix with elements t_{nk} .

- Again, the Laplace approximation can be used to optimize the hyperparameters in which the model and its Hessian are found using IRLS.

RVM For Classification

- ❑ The principal disadvantage is that the Hessian matrix has size $MK \times MK$, where M is the number of active basis functions, which gives an additional factor of K^3 in the computational cost of training compared with the two-class RVM.
- ❑ The principal disadvantage of the relevance vector machine is the relatively long training times compared with the SVM.
 - ✓ This is offset, however, by the avoidance of cross-validation runs to set the model complexity parameters.
 - ✓ Furthermore, because it yields sparser models, the computation time on test points, which is usually the more important consideration in practice, is typically much less.

References

- Scholkopf, B., A. Smola, R. C. Williamson, and P. L. Bartlett (2000). [New support vector algorithms](#). *Neural Computation* **12**(5), 1207–1245.
- Valiant, L. G. (1984). [A theory of the learnable](#). *Communications of the Association for Computing Machinery* **27**, 1134–1142.
- Anthony, M. and N. Biggs (1992). *An Introduction to Computational Learning Theory*. Cambridge University Press.
- Kearns, M. J. and U. V. Vazirani (1994). *An Introduction to Computational Learning Theory*. MIT Press.
- Vapnik, V. N. (1995). *The nature of statistical learning theory*. Springer.
- Vapnik, V. N. (1998). *Statistical learning theory*. Wiley.
- McAllester, D. A. (2003). [PAC-Bayesian stochastic model selection](#). *Machine Learning* **51**(1), 5–21.

References

- T. Fletcher, [Relevance Vector Machines Explained](#), October 2010.
- Tipping, M. (2001). [Sparse Bayesian learning and the relevance vector machine](#). *J. of Machine Learning Research* 1, 211–244.
- MacKay, D. (1995b). [Probable networks and plausible predictions — a review of practical Bayesian methods for supervised neural networks](#). *Network*.
- Neal, R. (1996). [Bayesian learning for neural networks](#). Springer.
- Bishop, C. and M. Tipping (2000). [Variational relevance vector machines](#). In *UAI*.
- Buntine, W. and A. Weigend (1991). [Bayesian backpropagation](#). *Complex Systems* 5, 603–643
- Tipping, M. (2001). [Sparse Bayesian learning and the relevance vector machine](#). *J. of Machine Learning Research* 1, 211–244.
- Wipf, D. and S. Nagarajan (2010, April). [Iterative Reweighted \$\ell_1\$ and \$\ell_2\$ Methods for Finding Sparse Solutions](#). *J. of Selected Topics in Signal Processing (Special Issue on Compressive Sensing)* 4(2).
- Wipf, D. and S. Nagarajan (2010, April). [Iterative Reweighted \$\ell_1\$ and \$\ell_2\$ Methods for Finding Sparse Solutions](#). *J. of Selected Topics in Signal Processing (Special Issue on Compressive Sensing)* 4(2).
- Wipf, D. and S. Nagarajan (2007). [A new view of automatic relevancy determination](#). In *NIPS*.
- MacKay, D. (1999). [Comparison of approximate methods for handling hyperparameters](#). *Neural Computation* 11(5), 1035–1068.
- Wipf, D. and S. Nagarajan (2007). [A new view of automatic relevancy determination](#). In *NIPS*.
- Wipf, D. and S. Nagarajan (2010, April). [Iterative Reweighted \$\ell_1\$ and \$\ell_2\$ Methods for Finding Sparse Solutions](#). *J. of Selected Topics in Signal Processing (Special Issue on Compressive Sensing)* 4(2).

References

- Wipf, D. and S. Nagarajan (2007). [A new view of automatic relevancy determination](#). In *NIPS*.
- Rasmussen, C. E. and J. Quiñonero-Candela (2005). [Healing the relevance vector machine by augmentation](#). In L. D. Raedt and S. Wrobel (Eds.), *Proceedings of the 22nd International Conference on Machine Learning*, pp. 689–696.
- T. Fletcher, [Relevance Vector Machines Explained](#), October 2010.
- Faul, A. C. and M. E. Tipping (2002). [Analysis of sparse Bayesian learning](#). In T. G. Dietterich, S. Becker, and Z. Ghahramani (Eds.), *Advances in Neural Information Processing Systems*, Volume 14, pp. 383–389. MIT Press.
- Tipping, M. E. and A. Faul (2003). [Fast marginal likelihood maximization for sparse Bayesian models](#). In C. M. Bishop and B. Frey (Eds.), *Proceedings Ninth International Workshop on Artificial Intelligence and Statistics*, Key West, Florida.
- Williams, O., A. Blake, and R. Cipolla (2005). [Sparse Bayesian learning for efficient visual tracking](#). *IEEE Transactions on Pattern Analysis and Machine Intelligence* **27**(8), 1292–1304.